

# On Optimal $k$ -Deletion Correcting Codes

Jin Sima and Jehoshua Bruck, *Fellow, IEEE*

**Abstract**—Levenshtein introduced the problem of constructing  $k$ -deletion correcting codes in 1966, proved that the optimal redundancy of those codes is  $O(k \log N)$  for constant  $k$ , and proposed an optimal redundancy single-deletion correcting code (using the so-called VT construction). However, the problem of constructing optimal redundancy  $k$ -deletion correcting codes remained open. Our key contribution is a major step towards a complete solution to this longstanding open problem for constant  $k$ . We present a  $k$ -deletion correcting code that has redundancy  $8k \log N + o(\log N)$  when  $k = o(\sqrt{\log \log N})$  and encoding/decoding algorithms of complexity  $O(N^{2k+1})$ .

**Index Terms**—Deletion codes, Varshamov-Tenengoltz code.

## I. INTRODUCTION

A set of binary vectors of length  $N$  is a  $k$ -deletion correcting code (denoted by  $\mathcal{C}$ ) iff any two vectors in  $\mathcal{C}$  do not share a subsequence of length  $N - k$ . The problem of constructing a  $k$ -deletion correcting code was introduced by Levenshtein [1]. He proved that the optimal redundancy (defined as  $N - \log |\mathcal{C}|$ ) is  $O(k \log N)$  for constant  $k$ . Specifically, it is in the range  $k \log N + o(\log N)$  to  $2k \log N + o(\log N)$  when  $k$  is a constant. In general, the redundancy  $O(k \log N)$  is orderwise optimal when  $k \leq O(N^\epsilon)$  for some  $\epsilon < 1$ . The optimal redundancy becomes  $O(k \log(N/k))$  when  $k = O(N)$  and  $k$  is small, e.g.,  $k = n/4$  [9]. When  $k \geq n/2$ , a  $k$ -deletion correcting code has cardinality at most two. In addition, Levenshtein proposed the following optimal construction (the well-known Varshamov-Tenengolts (VT) code [2]):

$$\left\{ (c_1, \dots, c_N) : \sum_{i=1}^N i c_i \equiv 0 \pmod{(N+1)} \right\}, \quad (1)$$

that is capable of correcting a single deletion with redundancy not more than  $\log(N+1)$  [1]. The encoding/decoding complexity of VT codes is linear in  $N$ . Generalizing the VT construction to correct more than a single deletion was elusive for more than 50 years. In particular, the past approaches [4], [5], [6] result in asymptotic code rates that are bounded away from 1.

A recent breakthrough paper [7] proposed a  $k$ -deletion correcting code construction with  $O(k^2 \log k \log N)$  redundancy and  $O_k(N \log^4 N)^1$  encoding/decoding complexity, where  $k$

is a constant. For the case  $k = 2$  deletions, the redundancy was improved in [12], [13]. Specifically, the code in [13] has redundancy of  $7 \log N$  and linear encoding/decoding complexity. The work in [14] considered correction with high probability and proposed a  $k$ -deletion correcting code construction with redundancy  $(k+1)(2k+1) \log N + o(\log N)$  and encoding/decoding complexity  $O(N^{k+1}/\log^{k-1} N)$ . The result for this randomized coding setting was improved in [8], where redundancy  $O(k \log(N/k))$  and complexity  $\text{poly}(N, k)$  were achieved. However, finding a deterministic  $k$ -deletion correcting code construction that achieves the optimal order redundancy  $O(k \log N)$  remained elusive.

Our key contribution is a solution to this longstanding open problem for constant  $k$ : We present a code construction that achieves  $O(k \log N)$  redundancy when  $k = o(\sqrt{\log \log N})$  and  $O(N^{2k+1})$  encoding/decoding computational complexity. Note that the complexity is polynomial in  $N$  when  $k$  is a constant. The following theorem summarizes our main result. We note that in this paper, the optimality of the code is redundancy-wise rather than cardinality-wise, namely, the focus is on the asymptotic rather than exact size of the code. The problem of finding optimal cardinality  $k$ -deletion correcting code appears highly nontrivial even for  $k = 1$ .

**Theorem 1.** *Let  $k$  and  $n$  be two integers satisfying  $k = o(\sqrt{\log \log n})$ . For integer  $N = n + 8k \log n + o(\log n)$ , there exists an encoding function  $\mathcal{E} : \{0, 1\}^n \rightarrow \{0, 1\}^N$ , computed in  $O(n^{2k+1})$  time, and a decoding function  $\mathcal{D} : \{0, 1\}^{N-k} \rightarrow \{0, 1\}^n$ , computed in  $O(n^{k+1}) = O(N^{k+1})$  time, such that for any  $\mathbf{c} \in \{0, 1\}^n$  and subsequence  $\mathbf{d} \in \{0, 1\}^{N-k}$  of  $\mathcal{E}(\mathbf{c})$ , we have that  $\mathcal{D}(\mathbf{d}) = \mathbf{c}$ .*

Recently, an independent work [9] proposed a  $k$ -deletion correcting code with  $O(k \log N)$  redundancy and better complexity of  $\text{poly}(N, k)$ . In contrast to a redundancy of  $8k \log N$  in this paper, the redundancy in [9] only specified asymptotically, and it is estimated to be at least  $200k \log n$ . Moreover, the techniques in [9] and this paper are different.

Here are the key building blocks in our code construction: (i) *generalizing the VT construction* to  $k$  deletions by considering constrained sequences, (ii) separating the encoded vector to blocks and using *concatenated codes* and (iii) a novel strategy to *separate the vector to blocks by a single pattern*. The following gives a more specific description of these ideas.

In our previous work for 2-deletions codes [13], we *generalized the VT construction*. In particular, we proved that while the higher order parity checks  $\sum_{i=1}^N i^j c_i \pmod{(N^j+1)}$ ,  $j = 0, 1, \dots, t$  do not work in general, those parity checks work in the two deletions case when the sequences are constrained to have no consecutive 1's. In this paper we generalize this idea, specifically, the higher order parity checks can be used to correct  $k = t/2$  deletions in sequences that satisfy the

Manuscript received October 15, 2019; revised June 4, 2020; accepted September 8, 2020. This work was supported in part by NSF grants CCF-1717884 and CCF-1816965. This paper was presented in part at the 2019 IEEE International Symposium on Information Theory.

J. Sima and J. Bruck are with the Electrical Engineering Department, California Institute of Technology, Pasadena, CA, 91125 USA (e-mail: jsima@caltech.edu ;bruck@caltech.edu).

Copyright (c) 2017 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

<sup>1</sup>The notion  $O_k$  denotes *parameterized complexity*, i.e.,  $O_k(N \log^4 N) = f(k)O(N \log^4 N)$  for some function  $f$ .



## II. OUTLINE AND PRELIMINARIES

In this section, we outline the ingredients that constitute our code construction and present notations that will be used throughout the paper. A summary of these notations is provided in Table I. We begin with an overview of the code construction, which has a concatenated code structure as described in Section I. In our construction, each codeword  $\mathbf{c}$  is splitted into blocks, with the boundaries between adjacent blocks given by synchronization patterns. Specifically, let  $t_1, \dots, t_J$  be the indices of the synchronization patterns in  $\mathbf{c}$ , i.e., the indices of the 1 entries in  $\mathbb{1}_{\text{sync}}(\mathbf{c})$ . Then the blocks are given by  $(c_{t_{j-1}+1}, \dots, c_{t_j-1})$  for  $j \in [0, J+1]$ , where  $t_j = 0$  if  $j = 0$  and  $t_j = n+1$  if  $j = J+1$ . The key idea of our construction is to use the VT generalization  $f$  (see Eq. (4) for definition) as parity checks to protect the synchronization vector  $\mathbb{1}_{\mathbf{c}}$ , and thus identify the indices of block boundaries. The proof details will be given in Lemma 1. Note that the parity  $f(\mathbf{c})$  in (4) has size  $O(k^2 \log n)$ . To compress the size of  $f(\mathbf{c})$  to the targeted  $O(k \log n)$ , in Lemma 2 we apply a modulo operation on the function  $f$ . Given the block boundaries, Lemma 3 provides an algorithm

TABLE I  
SUMMARY OF NOTATIONS

$\mathcal{B}_k(\mathbf{c})$	$\triangleq$	The set of sequences that share a length $n - k$ subsequence with $\mathbf{c} \in \{0, 1\}^n$ .
$\mathcal{R}_m$	$\triangleq$	The set of sequences where any two 1 entries are separated by at least $m - 1$ zeros.
$\mathbb{1}_{\text{sync}}(\mathbf{c})$	$\triangleq$	Synchronization vector defined in (2).
$\mathbf{m}^{(e)}$	$\triangleq$	Weights of order $e$ in the VT generalization, defined in (3).
$f(\mathbf{c})$	$\triangleq$	Generalized VT redundancy defined in (4).
$L$	$\triangleq$	Maximal length of zero runs in the synchronization vector of a $k$ -dense sequence, defined in (5)
$p(\mathbf{c})$	$\triangleq$	The function used to compute the redundancy protecting the synchronization vector $\mathbb{1}_{\text{sync}}(\mathbf{c})$ , defined in Lemma 2.
$\text{Hash}_k(\mathbf{c})$	$\triangleq$	The deletion correcting hash function for $k$ -dense sequences, defined in Lemma 3.
$T(\mathbf{c})$	$\triangleq$	The function that generates $k$ -dense sequences, defined in Lemma 4.
$H(\mathbf{c})$	$\triangleq$	Deletion correcting hash function for any sequence, defined in Lemma 5.

to protect the codewords. Specifically, we show that given the synchronization vector  $\mathbb{1}_{\text{sync}}(\mathbf{c})$ , it is possible to recover most of the blocks in  $\mathbf{c}$ , with up to  $2k$  block errors. These block errors can be corrected with the help of their deletion correcting hashes, which can be exhaustively computed and will be presented in Lemma 5. Hence, to correct the block errors, it suffices to protect the sequence of deletion correcting hashes using, for example, Reed-Solomon codes.

Lemma 2 and Lemma 3 together define a hash function that protects  $\mathbf{c}$  from  $k$  deletions. However, in order to exhaustively compute the block hash given in Lemma 5 and obtain the desired size of the redundancy, the block length has to be of order  $O(\text{poly}(k) \log n)$ . Hence the synchronization pattern must appear frequently enough in  $\mathbf{c}$ . Such sequence  $\mathbf{c}$  will be defined in the following as a  $k$ -dense sequence. To encode for

any given sequence  $\mathbf{c} \in \{0, 1\}^n$ , in Lemma 4 we present an invertible mapping that takes any sequence  $\mathbf{c} \in \{0, 1\}^n$  as input and outputs a  $k$ -dense sequence. Then, any binary sequence can be encoded into a  $k$ -dense sequence and protected.

Finally, the hash function defined by Lemma 2 and Lemma 3 is subject to deletion errors and has to be protected. To this end, we use an additional hash that encodes the deletion correcting hash of the hash function defined in Lemma 2 and Lemma 3. The additional hash is protected by a  $(k+1)$ -fold repetition code. Similar additional hash technique was also given in [7]. The final encoding/decoding algorithm of our code, which combines the results in Lemma 2, Lemma 3, and Lemma 4, is provided in Section III. The encoding is illustrated in Fig. 1.

Before presenting the lemmas, we give necessary definitions and notations. For a sequence  $\mathbf{c} \in \{0, 1\}^n$ , define its deletion ball  $\mathcal{B}_k(\mathbf{c})$  as the collection of sequences that share a length  $n - k$  subsequence with  $\mathbf{c}$ .

**Definition 2.** A sequence  $\mathbf{c} \in \{0, 1\}^n$  is said to be  $k$ -dense if the lengths of the 0 runs in  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  is at most

$$L \triangleq (\lceil \log k \rceil + 5)2^{\lceil \log k \rceil + 9} \lceil \log n \rceil + (3k + \lceil \log k \rceil + 4)(\lceil \log n \rceil + 9 + \lceil \log k \rceil). \quad (5)$$

For  $k$ -dense  $\mathbf{c}$ , the index distance between any two 1 entries in  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  is at most  $L + 1$ , i.e., the 0-runs between two 1 entries have length at most  $L + 1$ .

Note that  $f(\mathbf{c})$  is an integer vector that can be presented by  $\log(3k)^{6k+1}n^{(3k+1)(6k+1)}$  bits or by an integer in the range  $[0, (3k)^{6k+1}n^{(3k+1)(6k+1)} - 1]$ . In this paper, we interchangeably use  $f(\mathbf{c})$  to denote its binary presentation or integer presentation.

The following lemma shows that the synchronization vector  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  can be recovered from  $k$  deletions with the help of  $f(\mathbb{1}_{\text{sync}}(\mathbf{c}))$ . Its proof will be given in Section IV-A.

**Lemma 1.** For integers  $n$  and  $k$  and sequences  $\mathbf{c}, \mathbf{c}'$ , if  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$  and  $f(\mathbb{1}_{\text{sync}}(\mathbf{c})) = f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))$ , then  $\mathbb{1}_{\text{sync}}(\mathbf{c}) = \mathbb{1}_{\text{sync}}(\mathbf{c}')$ .

By virtue of Lemma 1, the synchronization vector  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  can be recovered from  $k$  deletions in  $\mathbf{c}$ , with the help of a hash  $f(\mathbb{1}_{\text{sync}}(\mathbf{c}))$  of size  $O(k^2 \log n)$  bits. To further reduce the size of the hash to  $O(k \log n)$  bits, we apply modulo operations on  $f(\mathbb{1}_{\text{sync}}(\mathbf{c}))$  in the following lemma, the proof of which will be proved in Section IV-B.

**Lemma 2.** For integers  $n$  and  $k = o(\sqrt{\log \log n})$ , there exists a function  $p : \{0, 1\}^n \rightarrow [1, 2^{2k \log n + o(\log n)}]$ , such that if  $f(\mathbb{1}_{\text{sync}}(\mathbf{c})) \equiv f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) \pmod{p(\mathbf{c})}$  for two sequences  $\mathbf{c} \in \{0, 1\}^n$  and  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$ , then  $\mathbb{1}_{\text{sync}}(\mathbf{c}) = \mathbb{1}_{\text{sync}}(\mathbf{c}')$ . Hence if

$$\begin{aligned} & (f(\mathbb{1}_{\text{sync}}(\mathbf{c})) \pmod{p(\mathbf{c})}, p(\mathbf{c})) \\ &= (f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) \pmod{p(\mathbf{c}'), p(\mathbf{c}')} \end{aligned}$$

and  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$ , we have that  $\mathbb{1}_{\text{sync}}(\mathbf{c}) = \mathbb{1}_{\text{sync}}(\mathbf{c}')$ .

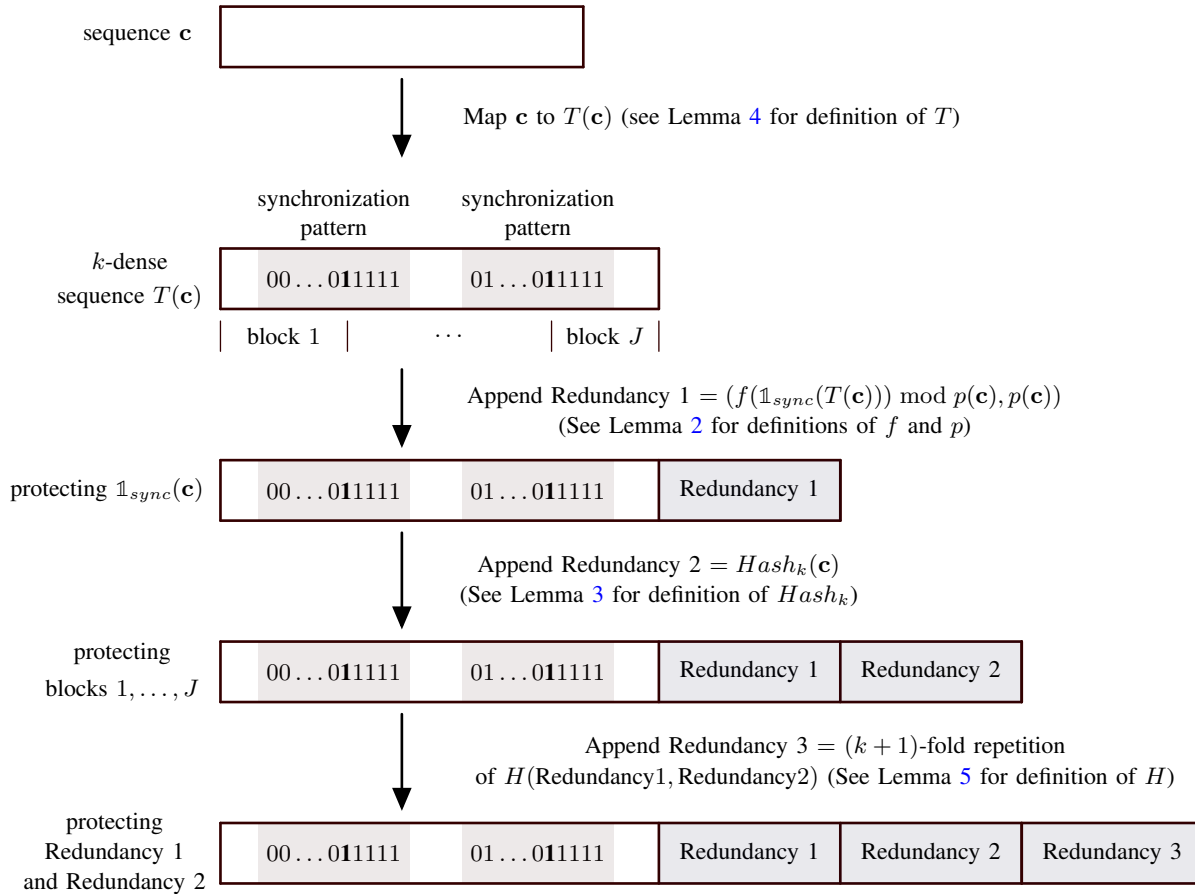


Fig. 1. Illustrating the encoding procedure of our  $k$ -deletion correcting code construction, based on results of Lemma 2, Lemma 3, Lemma 4, and Lemma 5. The indices of the bold 1 entries are the indices of the 1 entries in the synchronization vector  $\mathbb{1}_{sync}(T(\mathbf{c}))$ .

Lemma 2 presents a hash of size  $4k \log n + o(\log n)$  bits for correcting  $\mathbb{1}_{sync}(\mathbf{c})$ . With the knowledge of the synchronization vector  $\mathbb{1}_{sync}(\mathbf{c})$ , the next lemma shows that the sequence  $\mathbf{c}$  can be further recovered using another  $4k \log n + o(\log n)$  bit hash, when  $\mathbf{c}$  is  $k$ -dense, i.e., when the synchronization pattern occurs frequently enough in  $\mathbf{c}$ . The proof of Lemma 3 will be given in Section V.

**Lemma 3.** *For integers  $n$  and  $k = o(\sqrt{\log \log n})$ , there exists a function  $Hash_k : \{0, 1\}^n \rightarrow \{0, 1\}^{4k \log n + o(\log n)}$ , such that every  $k$ -dense sequence  $\mathbf{c} \in \{0, 1\}^n$  can be recovered, given its synchronization vector  $\mathbb{1}_{sync}(\mathbf{c})$ , its length  $n - k$  subsequence  $\mathbf{d}$ , and  $Hash_k(\mathbf{c})$ .*

Combining Lemma 2 and Lemma 3, we obtain size  $O(k \log n)$  hash function to correct deletions for a  $k$ -dense sequence. To encode for arbitrary sequence  $\mathbf{c} \in \{0, 1\}^n$ , a mapping that transforms any sequence to a  $k$ -dense sequence is given in the following lemma. The details will be given in Section VI-B.

**Lemma 4.** *For integers  $k$  and  $n > k$ , there exists a map  $T : \{0, 1\}^n \rightarrow \{0, 1\}^{n+3k+3\lceil \log k \rceil + 15}$ , computable in  $\text{poly}(n, k)$  time, such that  $T(\mathbf{c})$  is a  $k$ -dense sequence for  $\mathbf{c} \in \{0, 1\}^n$ . Moreover, the sequence  $\mathbf{c}$  can be recovered from  $T(\mathbf{c})$ .*

The next three lemmas (Lemma 5, Lemma 6, and Lemma 7)

present existence results that are necessary to prove Lemma 2 and Lemma 3. Lemma 5 gives a  $k$ -deletion correcting hash function for short sequences, which is the block hash described above in proving Lemma 3. It is an extension of the result in [7]. Lemma 6 is a slight variation of the result in [1]. It shows the equivalence between correcting deletions and correcting deletions and insertions. Lemma 6 will be used in the proof of Lemma 1, where we need an upper bound on the number of deletions/insertions in  $\mathbb{1}_{sync}(\mathbf{c})$  caused by a deletion in  $\mathbf{c}$ . Lemma 7 (see [15]) gives an upper bound on the number of divisors of a positive integer  $n$ . With Lemma 7, we show that the VT generalization in Lemma 1 can be compressed by taking modulo operations. The details will be given in the proof of Lemma 2.

**Lemma 5.** *For any integers  $w$ ,  $n$ , and  $k$ , there exists a hash function  $H : \{0, 1\}^w \rightarrow \{0, 1\}^{\lceil (w/\lceil \log n \rceil) \rceil (2k \log \log n + O(1))}$ , computable in  $O_k((w/\log n)n \log^{2k} n)$  time, such that any sequence  $\mathbf{c} \in \{0, 1\}^w$  can be recovered from its length  $w - k$  subsequence  $\mathbf{d}$  and the hash  $H(\mathbf{c})$ .*

*Proof.* We first show by counting arguments the existence of a hash function  $H' : \{0, 1\}^{\lceil \log n \rceil} \rightarrow \{0, 1\}^{2k \log \log n + O(1)}$ , exhaustively computable in  $O_k(n \log^{2k} n)$  time, such that



$H'(s) \neq H'(s')$  for all  $s \in \{0, 1\}^{\lceil \log n \rceil}$  and  $s' \in \mathcal{B}_k(s) \setminus \{s\}$ . The hash  $H'(c')$  protects the sequence  $s \in \{0, 1\}^{\lceil \log n \rceil}$  from  $k$  deletions. Note that  $|\mathcal{B}_k(c')| \leq \binom{\lceil \log n \rceil}{k} 2^k \leq 2 \lceil \log n \rceil^{2k}$ . Hence it suffices to use brute force and greedily assign a hash value for each sequence  $s \in \{0, 1\}^{\lceil \log n \rceil}$  such that  $H'(s) \neq H'(s')$  for all  $s' \in \mathcal{B}_k(s) \setminus \{s\}$ . Since the size of  $\mathcal{B}_k(s)$  is upper bounded by  $2 \lceil \log n \rceil^{2k}$ , there always exists such a hash  $H'(s) \in \{0, 1\}^{\log(2 \lceil \log n \rceil^{2k+1})}$ , that has no conflict with the hash values of sequences in  $\mathcal{B}_k(s)$ . The total complexity is  $O_k(n \log^{2k} n)$  and the size of the hash value  $H'(s)$  is  $2k \log \log n + O(1)$ .

Now split  $c$  into  $\lceil (w/\lceil \log n \rceil) \rceil$  blocks  $c_{(i-1)\lceil \log n \rceil+1}, \dots, c_{i\lceil \log n \rceil}$ ,  $i \in [1, \lceil (w/\lceil \log n \rceil) \rceil]$  of length  $\lceil \log n \rceil$ . If the length of the last block is less than  $\lceil \log n \rceil$ , add zeros to the end of the last block such that its length is  $\lceil \log n \rceil$ . Assign a hash value  $\mathbf{h}_i = H'((c_{(i-1)\lceil \log n \rceil+1}, \dots, c_{i\lceil \log n \rceil}))$ ,  $i \in [1, \lceil (w/\lceil \log n \rceil) \rceil]$  for each block. Let  $H(c) = (\mathbf{h}_1, \dots, \mathbf{h}_{\lceil (w/\lceil \log n \rceil) \rceil})$  be the concatenation of  $\mathbf{h}_i$  for  $i \in [1, \lceil (w/\lceil \log n \rceil) \rceil]$ .

We show that  $H(c)$  protects  $c$  from  $k$  deletions. Let  $d$  be a length  $n - k$  subsequence of  $c$ . Note that  $d_{(i-1)\lceil \log n \rceil+1}$  and  $d_{i\lceil \log n \rceil-k}$  come from bits  $c_{(i-1)\lceil \log n \rceil+1+x}$  and  $c_{i\lceil \log n \rceil-k+y}$  respectively after deletions in  $c$ , where the integers  $x, y \in [0, k]$ . Therefore,  $(d_{(i-1)\lceil \log n \rceil+1}, \dots, d_{i\lceil \log n \rceil-k})$  is a length  $\lceil \log n \rceil - k$  subsequence of  $(c_{(i-1)\lceil \log n \rceil+1+x}, \dots, c_{i\lceil \log n \rceil-k+y})$ , and thus a subsequence of the block  $(c_{(i-1)\lceil \log n \rceil+1}, \dots, c_{i\lceil \log n \rceil})$ . Hence the  $i$ -th block  $(c_{(i-1)\lceil \log n \rceil+1}, \dots, c_{i\lceil \log n \rceil})$  can be recovered from  $\mathbf{h}_i = H'(c_{(i-1)\lceil \log n \rceil+1}, \dots, c_{i\lceil \log n \rceil})$  and  $(d_{(i-1)\lceil \log n \rceil+1}, \dots, d_{i\lceil \log n \rceil-k})$ . Therefore,  $c$  can be recovered given  $d$  and  $H(c)$ . The length of  $H(c)$  is  $\lceil (w/\lceil \log n \rceil) \rceil (2k \log \log n + O(1))$  and the complexity of  $H(c)$  is  $O_k((w/\log n)n \log^{2k} n)$ .  $\square$

**Lemma 6.** Let  $r, s$ , and  $k$  be integers satisfying  $r + s \leq k$ . For sequences  $c, c' \in \{0, 1\}^n$ , if  $c'$  and  $c$  share a common resulting sequence after  $r$  deletions and  $s$  insertions in both, then  $c' \in \mathcal{B}_k(c)$ .

**Lemma 7.** For a positive integer  $n \geq 3$ , the number of divisors of  $n$  is upper bounded by  $2^{1.6 \ln n / (\ln \ln n)}$ .

The proofs of Lemma 1, Lemma 2, Lemma 3, and Lemma 4 rely on several propositions, the details of which will be presented in the next sections. For convenience, a dependency graph for the theorem, lemmas and propositions is given in Fig. 2.

### III. PROOF OF THEOREM 1

Based on the lemmas stated in Section II, in this section we present the encoding function  $\mathcal{E}$  and the decoding function  $\mathcal{D}$  of our  $k$ -deletion correcting code, and prove Theorem 1. Given any sequence  $c \in \{0, 1\}^n$ , let the function  $\mathcal{E} : \{0, 1\}^n \rightarrow \{0, 1\}^{n+8k \log n + o(\log n)}$  be given by

$$\mathcal{E}(c) = (T(c), R'(c), R''(c)),$$

where

$$R'(c) = (f(\mathbb{1}_{\text{sync}}(T(c))) \bmod p(T(c))),$$

$$p(T(c)), \text{Hash}_k(T(c)), \text{ and } R''(c) = \text{Rep}_{k+1}(H(R'(c))).$$

Function  $\text{Rep}_{k+1}(H(R'(c)))$  is the  $(k+1)$ -fold repetition of the bits in  $H(R'(c))$ , where function  $H(R'(c))$  is defined in Lemma 5 and protects  $R'(c)$  from  $k$  deletions. Function  $T(c)$  is defined in Lemma 4 and transforms  $c$  into a  $k$ -dense sequence (see Definition 2 for definition of a  $k$ -dense sequence). Function  $f(\mathbb{1}_{\text{sync}}(T(c)))$  in  $R'(c)$  is represented by an integer and protects the synchronization vector  $\mathbb{1}_{\text{sync}}(c)$  by Lemma 1. Function  $p(T(c))$  is defined in Lemma 2, which compresses the hash  $f(\mathbb{1}_{\text{sync}}(T(c)))$ . Function  $\text{Hash}_k(T(c))$  is defined in Lemma 3 and protects a  $k$ -dense sequence from  $k$  deletions.

Note that  $k = o(\sqrt{\log \log n})$ . Hence, according to Lemma 2, Lemma 3, and Lemma 4, the length of  $R'(c)$  is  $N_1 = 8k \log(n + 3k + 3 \lceil \log k \rceil + 15) + o(\log(n + 3k + 3 \lceil \log k \rceil + 15)) = 8k \log n + o(\log n)$ . The length of  $R''(c)$  is  $N_2 = 2k(k+1)(N_1/\lceil \log n \rceil) \log \log n = o(\log n)$ . The length of  $T(c)$  is  $n + N_0 = n + 3k + 3 \lceil \log k \rceil + 15$ . Therefore, the length of  $\mathcal{E}(c)$  is  $n + N_0 + N_1 + N_2 = n + 8k \log n + o(\log n)$ . The redundancy of the code is  $8k \log n + o(\log n)$ .

To show how the sequence  $c$  can be recovered from a length  $N - k$  subsequence  $d$  of  $\mathcal{E}(c)$  and implement the computation of the decoding function  $\mathcal{D}(d)$ , we prove that

- 1) **Statement 1:** The redundancy  $R'(c)$  can be recovered given  $(d_{n+N_0+N_1+1}, \dots, d_{n+N_0+N_1+N_2-k})$ .
- 2) **Statement 2:** The sequence  $c$  can be recovered given  $(d_1, \dots, d_{n+N_0-k})$  and  $R'(c)$ .

We first prove Statement 1. Note that  $d_{n+N_0+N_1+1}$  and  $d_{n+N_0+N_1+N_2-k}$  come from bits  $\mathcal{E}(c)_{n+N_0+N_1+1+x}$  and  $\mathcal{E}(c)_{n+N_0+N_1+N_2-k+y}$  respectively after deletions in  $\mathcal{E}(c)$ , where  $x, y \in [0, k]$ . Hence  $(d_{n+N_0+N_1+1}, \dots, d_{n+N_0+N_1+N_2-k})$  is a length  $N_2 - k$  subsequence of  $(\mathcal{E}(c)_{n+N_0+N_1+1+x}, \dots, \mathcal{E}(c)_{n+N_0+N_1+N_2-k+y})$ , and thus a subsequence of  $(\mathcal{E}(c)_{n+N_0+N_1+1}, \dots, \mathcal{E}(c)_{n+N_0+N_1+N_2}) = R''(c)$ . Since  $R''(c)$  is  $k+1$ -fold repetition code and thus a  $k$ -deletion correcting code that protects  $H(R'(c))$ , the hash function  $H(R'(c))$  can be recovered from  $(d_{n+N_0+N_1+1}, \dots, d_{n+N_0+N_1+N_2-k})$ .

Similarly,  $(d_{n+N_0+1}, \dots, d_{n+N_0+N_1-k})$  is a length  $N_1 - k$  subsequence of  $(\mathcal{E}(c)_{n+N_0+1}, \dots, \mathcal{E}(c)_{n+N_0+N_1}) = R'(c)$ . From Lemma 5, the function  $R'(c)$  can be recovered from  $H(R'(c))$  and  $(d_{n+N_0+1}, \dots, d_{n+N_0+N_1-k})$ . Hence Statement 1 holds.

We now prove Statement 2. Note that  $R'(c)$  contains hashes  $(f(\mathbb{1}_{\text{sync}}(T(c))) \bmod p(T(c)), p(T(c)))$  and  $\text{Hash}_k(T(c))$ . Moreover,  $(d_1, \dots, d_{n+N_0-k})$  is a length  $n + N_0 - k$  subsequence of  $(\mathcal{E}(c)_1, \dots, \mathcal{E}(c)_{n+N_0}) = T(c)$ . According to Lemma 2, the synchronization vector  $\mathbb{1}_{\text{sync}}(T(c))$  can be recovered from hash  $(f(\mathbb{1}_{\text{sync}}(T(c))) \bmod p(T(c)), p(T(c)))$  and  $(d_1, \dots, d_{n+N_0-k})$ , by exhaustively searching over all length  $n + N_0$  supersequence  $c'$  of  $(d_1, \dots, d_{n+N_0-k})$  such that  $f(\mathbb{1}_{\text{sync}}(c')) = f(\mathbb{1}_{\text{sync}}(T(c)))$ . Then we have that  $\mathbb{1}_{\text{sync}}(T(c)) = \mathbb{1} + \text{sync}(c')$ . Since by Lemma 4,  $T(c)$  is a  $k$ -dense sequence, by Lemma 3, it can be recovered from  $\mathbb{1}_{\text{sync}}(T(c))$ ,  $\text{Hash}_k(T(c))$ , and the length  $n + N_0 - k$  subsequence  $(d_1, \dots, d_{n+N_0-k})$  of  $T(c)$ . Finally, the

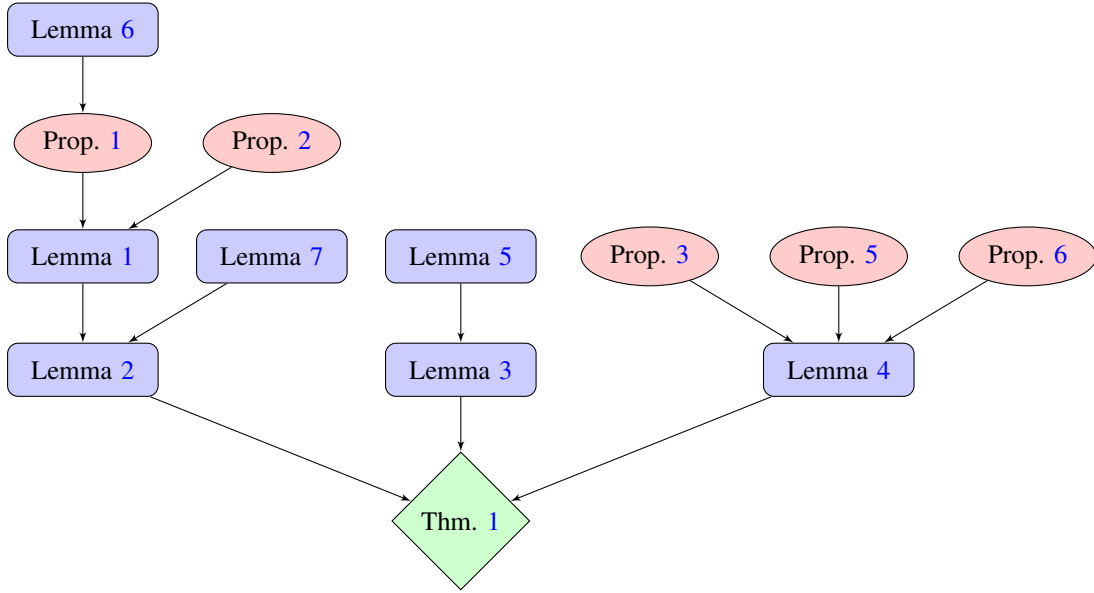


Fig. 2. Dependencies of the claims in the paper.

sequence  $\mathbf{c}$  can be recovered from  $T(\mathbf{c})$  by Lemma 4. Hence Statement 2 holds and  $\mathbf{c}$  can be recovered.

The encoding complexity of  $\mathcal{E}(\mathbf{c})$  is  $O(n^{2k+1})$ , which comes from brute force search for integer  $p(T(\mathbf{c}))$ . The decoding complexity is  $O(n^k + 1)$ , which comes from brute force search for the correct  $\mathbb{1}_{sync}(T(\mathbf{c}))$ , given  $f(\mathbb{1}_{sync}(T(\mathbf{c}))) \bmod p(T(\mathbf{c}))$  and  $p(T(\mathbf{c}))$ .

#### IV. PROTECTING THE SYNCHRONIZATION VECTORS

In this section we present a hash function with size  $4k \log n + o(\log n)$  to protect the synchronization vector  $\mathbb{1}_{sync}(\mathbf{c})$  from  $k$  deletions in  $\mathbf{c}$  and prove Lemma 2. We first prove Lemma 1, which is decomposed to Proposition 1 and Proposition 2. In Proposition 1 we present an upper bound on the radius of the deletion ball for the synchronization vector. In Proposition 2, we prove that the higher order parity check helps correct multiple deletions for sequences in which there is a 0-run of length at least  $3k - 1$  between any two 1's. Since  $\mathbb{1}_{sync}(\mathbf{c})$  is such a sequence, we conclude that the higher order parity check helps recover  $\mathbb{1}_{sync}(\mathbf{c})$ . After obtaining a bound on the difference between the higher order parity checks of two ambiguous sequence, we then apply Proposition 2 on the synchronization vector  $\mathbb{1}_{sync}(\mathbf{c})$  to prove Lemma 1, which replaces the higher order parity checks in Proposition 2 by the higher parity checks modulo a numbers. After proving Lemma 1, we use Lemma 7 to further compress the size of the higher order parity check that protects  $\mathbb{1}_{sync}(\mathbf{c})$  and then prove Lemma 2.

**Proposition 1.** For  $\mathbf{c}, \mathbf{c}' \in \{0, 1\}^n$ , if  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$ , then  $\mathbb{1}_{sync}(\mathbf{c}') \in \mathcal{B}_{3k}(\mathbb{1}_{sync}(\mathbf{c}))$ .

*Proof.* Since  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$ , the sequences  $\mathbf{c}'$  and  $\mathbf{c}$  share a common subsequence after  $k$  deletions in both. We now show that a single deletion in  $\mathbf{c}$  causes at most two deletions and one insertion in its synchronization vector  $\mathbb{1}_{sync}(\mathbf{c})$ . We first show that a deletion in  $\mathbf{c}$  can destroy and generate

at most 1 synchronization pattern. This is because for any synchronization pattern that is destroyed or generated, there must be a deletion that occurs within the synchronization pattern. Hence any two destroyed or generated synchronization patterns cannot be caused by the same deletion. Therefore, we need to consider four cases in total. Let  $\mathbf{d}'$  be the subsequence of  $\mathbf{c}$  after a single deletion.

- 1) The deletion destroys a synchronization pattern  $(c_{i+1}, \dots, c_{i+3k+\lceil \log k \rceil+4})$  for some  $i$  and no synchronization pattern is generated. Then the sequence  $\mathbb{1}_{sync}(\mathbf{d}')$  can be obtained by deleting the 1 entry  $\mathbb{1}_{sync}(\mathbf{c})_{i+3k}$  in  $\mathbb{1}_{sync}(\mathbf{c})$ .
- 2) The deletion generates a new synchronization pattern  $(c'_{i'+1}, \dots, c'_{i'+3k+\lceil \log k \rceil+4})$  for some  $i'$  and destroys a synchronization pattern  $(c_{i+1}, \dots, c_{i+3k+\lceil \log k \rceil+4})$ . The sequence  $\mathbb{1}_{sync}(\mathbf{d}')$  can be obtained by deleting the 1 entry  $\mathbb{1}_{sync}(\mathbf{c})_{i+3k}$  and the 0 entry  $\mathbb{1}_{sync}(\mathbf{c})_{i+3k-1}$  in  $\mathbb{1}_{sync}(\mathbf{c})$  and inserting a 1 entry at  $\mathbb{1}_{sync}(\mathbf{c})_{i'+3k}$ .
- 3) The deletion generates a new synchronization pattern  $(c'_{i'+1}, \dots, c'_{i'+3k+\lceil \log k \rceil+4})$  for some  $i'$  and no synchronization pattern is destroyed. Then the  $\mathbb{1}_{sync}(\mathbf{d}')$  can be obtained by deleting two 0 entries  $\mathbb{1}_{sync}(\mathbf{c})_{i'+3k}$  and  $\mathbb{1}_{sync}(\mathbf{c})_{i'+3k+1}$  in  $\mathbb{1}_{sync}(\mathbf{c})$  and inserting a 1 entry at  $\mathbb{1}_{sync}(\mathbf{c})_{i'+3k}$ .
- 4) No synchronization pattern is generated or destroyed. Then  $\mathbb{1}_{sync}(\mathbf{d}')$  can be obtained by deleting a 0 entry  $\mathbb{1}_{sync}(\mathbf{c})_j$ , where  $j$  is the location of the deletion.

In summary, in each of the above cases, a single deletion in  $\mathbf{c}$  causes at most two deletions and one insertion in  $\mathbb{1}_{sync}(\mathbf{c})$ . Hence  $k$  deletions in  $\mathbf{c}$  and  $\mathbf{c}'$  cause at most  $2k$  deletions and  $k$  insertions in  $\mathbb{1}_{sync}(\mathbf{c})$  and  $\mathbb{1}_{sync}(\mathbf{c}')$  respectively. According to Lemma 6, we have that  $\mathbb{1}_{sync}(\mathbf{c}') \in \mathcal{B}_{3k}(\mathbb{1}_{sync}(\mathbf{c}))$  when  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$ . Hence, Proposition 1 is proved.  $\square$

Let  $\mathcal{R}_m$  be the set of length  $n$  sequences in which there

is a 0 run of length at least  $m - 1$  between any two 1's. Any two 1's in a sequence  $\mathbf{c} \in \mathcal{R}_m$  have index distance at least  $m$ . The following lemma shows that the sequences in  $\mathcal{R}_{3k}$  can be protected using higher order parity checks. Note that compared to the higher order parity checks  $f(\mathbf{c})$ , the higher order parity checks in the following proposition do not have modulo operations.

**Proposition 2.** *For sequences  $\mathbf{c}, \mathbf{c}' \in \mathcal{R}_{3k}$ , if  $\mathbf{c}' \in \mathcal{B}_{3k}(\mathbf{c})$  and  $\mathbf{c} \cdot \mathbf{m}^{(e)} = \mathbf{c}' \cdot \mathbf{m}^{(e)}$  for  $e \in [0, 6k]$ , then  $\mathbf{c} = \mathbf{c}'$ .*

*Proof.* We first compute the difference  $\mathbf{c} \cdot \mathbf{m}^{(e)} - \mathbf{c}' \cdot \mathbf{m}^{(e)}$ ,  $e \in [0, 6k]$ . Since  $\mathbf{c}' \in \mathcal{B}_{3k}(\mathbf{c})$ , there exist two subsets  $\delta = \{\delta_1, \dots, \delta_{3k}\} \subset [1, n]$  and  $\delta' = \{\delta'_1, \dots, \delta'_{3k}\} \subset [1, n]$  such that deleting bits with indices  $\delta$  and  $\delta'$  respectively from  $\mathbf{c}$  and  $\mathbf{c}'$  results in the same length  $n - 3k$  subsequence, i.e.,  $(c_i : i \notin \delta) = (c'_i : i \notin \delta')$ . Let  $\Delta = \{i : c_i = 1\}$  and  $\Delta' = \{i : c'_i = 1\}$  be the indices of 1 entries in  $\mathbf{c}$  and  $\mathbf{c}'$  respectively. Let  $S_1 = \Delta \cap ([1, n] \setminus \delta)$  be the indices of 1 entries that are deleted in  $\mathbf{c}$ . Then  $S_1^c = \Delta \cap ([1, n] \setminus \delta)$  denotes the indices of 1 entries that are not deleted. Similarly, let  $S_2 = \Delta' \cap \delta'$  and  $S_2^c = \Delta' \cap ([1, n] \setminus \delta')$  be the indices of 1 entries that are deleted and not in  $\mathbf{c}'$  respectively. Let the elements in  $\delta \cup \delta'$  be ordered by  $1 \leq p_1 \leq p_2 \leq \dots \leq p_{6k} \leq n$ . Denote  $p_0 = 0$  and  $p_{6k+1} = n$ . Then we have that

$$\begin{aligned} & \mathbf{c} \cdot \mathbf{m}^{(e)} - \mathbf{c}' \cdot \mathbf{m}^{(e)} \\ &= \sum_{\ell \in \Delta} \mathbf{m}_\ell^{(e)} - \sum_{\ell \in \Delta'} \mathbf{m}_\ell^{(e)} \\ &= \sum_{\ell \in \Delta} \left( \sum_{i=1}^{\ell} i^e \right) - \sum_{\ell \in \Delta'} \left( \sum_{i=1}^{\ell} i^e \right) \\ &= \sum_{i=1}^n \left( \sum_{\ell \in \Delta \cap [i, n]} i^e \right) - \sum_{i=1}^n \left( \sum_{\ell \in \Delta' \cap [i, n]} i^e \right) \\ &= \sum_{i=1}^n (|\Delta \cap [i, n]| - |\Delta' \cap [i, n]|) i^e \\ &= \sum_{i=1}^n (|S_1 \cap [i, n]| + |S_1^c \cap [i, n]| - |S_2 \cap [i, n]| \\ &\quad - |S_2^c \cap [i, n]|) i^e \\ &= \sum_{j=0}^{6k} \sum_{i=p_j+1}^{p_{j+1}} (|S_1 \cap [i, n]| - |S_2 \cap [i, n]| + |S_1^c \cap [i, n]| \\ &\quad - |S_2^c \cap [i, n]|) i^e \\ &\stackrel{(a)}{=} \sum_{j=0}^{6k} \sum_{i=p_j+1}^{p_{j+1}} (|S_1 \cap [p_{j+1}, n]| \\ &\quad - |S_2 \cap [p_{j+1}, n]| + |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]|) i^e, \quad (6) \end{aligned}$$

where (a) holds since by definition of  $p_j$ , there is no deleted 1 entry in interval  $(p_j, p_{j+1}) = \{p_j + 1, \dots, p_{j+1} - 1\}$ ,  $j \in [0, 6k]$ . In the following we show

**Statement 1:**  $-1 \leq |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \leq 1$  for  $i \in [1, n]$ .

**Statement 2:** For each interval  $(p_j, p_{j+1}] = \{p_j + 1, \dots, p_{j+1}\}$ ,  $j = 0, \dots, 6k$ , we have either  $|S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \leq 0$  for all  $i \in (p_j, p_{j+1}]$  or  $|S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \geq 0$  for all  $i \in (p_j, p_{j+1}]$ .

We first prove **Statement 1**. Note that deleting bits with indices  $\delta$  in  $\mathbf{c}$  and deleting bits with indices  $\delta'$  in  $\mathbf{c}'$  result in the same subsequence. Hence, for every  $i \in S_1^c$ , there is a unique corresponding index  $i' \in S_2^c$  such that the two 1 entries  $c_i$  and  $c'_{i'}$  end in the same location after deletions, i.e.,  $i - |\delta \cap [1, i - 1]| = i' - |\delta' \cap [1, i' - 1]|$ . This implies that  $|i' - i| \leq 3k$ . Fix integers  $i$  and  $i'$ . Then by definition of  $i$  and  $i'$ , for every  $x \in S_1^c \cap [i + 1, n]$ , there is a unique corresponding  $y \in S_2^c \cap [i' + 1, n]$  such that the two 1 entries  $c_x$  and  $c'_y$  end in the same location after deletions. Therefore, we have that  $|S_1^c \cap [i + 1, n]| = |S_2^c \cap [i' + 1, n]|$ , and thus that  $|S_1^c \cap [i, n]| = |S_2^c \cap [i', n]|$ . If  $i' \geq i$ , then

$$\begin{aligned} & |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \\ &= |S_1^c \cap [i, n]| - |S_2^c \cap [i', n]| - |S_2^c \cap [i, i' - 1]| \\ &= -|S_2^c \cap [i, i' - 1]| \\ &\stackrel{(a)}{\geq} -|S_2^c \cap [i, i + 3k - 1]| \\ &\stackrel{(b)}{\geq} -1, \end{aligned}$$

where (a) follows from the fact that  $i' \leq i + 3k$  and (b) follows from the fact that  $\mathbf{c}, \mathbf{c}' \in \mathcal{R}_{3k}$ . Also we have that  $|S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| = -|S_2^c \cap [i, i' - 1]| \leq 0$ . Hence when  $i' \leq i$ , we have that  $-1 \leq |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \leq 0$ . Similarly, when  $i' < i$ , we have that

$$\begin{aligned} & |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \\ &= |S_1^c \cap [i, n]| - |S_2^c \cap [i', n]| + |S_2^c \cap [i', i - 1]| \\ &= |S_2^c \cap [i', i - 1]| \\ &\leq |S_2^c \cap [i', i' + 3k - 1]| \\ &\leq 1, \end{aligned}$$

and that  $|S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| = |S_2^c \cap [i', i - 1]| \geq 0$ . Therefore, we have that  $0 \leq |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \leq 1$  when  $i' < i$ . Thus **Statement 1** is proved.

We now prove **Statement 2** by contradiction. Suppose on the contrary, there exist  $i_1, i_2 \in (p_j, p_{j+1}]$  such that  $i_1 < i_2$  and

$$(|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]|)(|S_1^c \cap [i_2, n]| - |S_2^c \cap [i_2, n]|) < 0$$

From **Statement 1** we have that  $|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]| \in [-1, 1]$  and that  $|S_1^c \cap [i_2, n]| - |S_2^c \cap [i_2, n]| \in [-1, 1]$ . Hence by symmetry it can be assumed that  $|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]| = -1$  and  $|S_1^c \cap [i_2, n]| - |S_2^c \cap [i_2, n]| = 1$ . As shown in proof of **Statement 1**, for every element  $i \in S_1^c$ , there is a corresponding element  $i' \in S_2^c$  such that the two 1 entries  $c_i$  and  $c'_{i'}$  end in the same location after deletions. Hence, for  $y = \min_{i \in S_2^c \cap [i_1, n]} i$ , there exists an integer  $x \in S_1^c$  such that the two 1 entries  $c_x$  and  $c'_y$  are in the same location after deletions, i.e.,  $x - |\delta \cap [1, x - 1]| = y - |\delta' \cap [1, y - 1]|$ . Since  $|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]| = -1$ , we have that  $x \in S_1^c \cap [1, i_1 - 1]$ . Otherwise, we have that  $x \in S_1^c \cap [i_1, n]$  and for every integer  $i' \in S_2^c \cap (y, n]$ , there exists an integer  $i \in S_1^c \cap (x, n]$  such that  $c_i$  and  $c'_{i'}$  end up in the same location after deletions. This implies that  $|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]| \geq 0$ ,

contradicting the fact that  $|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]| = -1$ . Therefore,

$$\begin{aligned} i_1 - |\delta \cap [1, i_1 - 1]| &> i_1 - 1 - |\delta \cap [1, i_1 - 1]| \\ &\geq x - |\delta \cap [1, x - 1]| \\ &= y - |\delta' \cap [1, y - 1]| \\ &\geq i_1 - |\delta' \cap [1, i_1 - 1]|, \end{aligned}$$

which implies that

$$|\delta \cap [1, i_1 - 1]| < |\delta' \cap [1, i_1 - 1]|. \quad (7)$$

Similarly, from  $|S_1^c \cap [i_2, n]| - |S_2^c \cap [i_2, n]| = 1$  we have that

$$|\delta \cap [1, i_2 - 1]| > |\delta' \cap [1, i_2 - 1]|. \quad (8)$$

Eq. (7) and Eq. (8) implies that

$$\begin{aligned} &|\delta \cap [1, i_2 - 1]| - |\delta \cap [1, i_1 - 1]| \\ &\geq |\delta' \cap [1, i_2 - 1]| + 1 - |\delta' \cap [1, i_1 - 1]| + 1 \\ &\geq 2. \end{aligned} \quad (9)$$

However, since  $i_1, i_2 \in (p_j, p_{j+1}]$  and no deletion occurs in the interval  $(p_j, p_{j+1}]$ , we have that  $|\delta \cap [1, i_1]| = |\delta \cap [1, i_2 - 1]|$  and  $|\delta' \cap [1, i_1]| = |\delta' \cap [1, i_2 - 1]|$ , which implies that

$$\begin{aligned} &|\delta \cap [1, i_2 - 1]| - |\delta \cap [1, i_1 - 1]| \\ &\leq |\delta \cap [1, i_2 - 1]| - |\delta \cap [1, i_1]| + 1 \\ &= 1, \end{aligned}$$

contradicting Eq. (9). Hence there do not exist different integers  $i_1, i_2 \in (p_j, p_{j+1}]$  such that

$$(|S_1^c \cap [i_1, n]| - |S_2^c \cap [i_1, n]|)(|S_1^c \cap [i_2, n]| - |S_2^c \cap [i_2, n]|) < 0.$$

Hence **Statement 2** is proved.

Now we continue to prove Proposition 2. Denote

$$s_i \triangleq |S_1 \cap [i, n]| - |S_2 \cap [i, n]| + |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]|. \quad (10)$$

Note that no deletion occurs in the interval  $(p_j, p_{j+1}]$ , it follows that

$$\begin{aligned} &|S_1 \cap [i, n]| - |S_2 \cap [i, n]| \\ &= |S_1 \cap [p_{j+1}, n]| - |S_2 \cap [p_{j+1}, n]| \end{aligned} \quad (11)$$

for  $i \in (p_j, p_{j+1}]$ . Combining (11) with **Statement 1** and **Statement 2**, we conclude that for each interval  $(p_j, p_{j+1}]$ ,  $j \in \{0, \dots, 6k\}$ , either  $s_i \geq 0$  for all  $i \in (p_j, p_{j+1}]$  or  $s_i \leq 0$  for all  $i \in (p_j, p_{j+1}]$ . Let  $\mathbf{x} = (x_0, \dots, x_{6k}) \in \{-1, 1\}^{6k+1}$  be a vector defined by

$$x_i = \begin{cases} -1, & \text{if } s_j < 0 \text{ for some } j \in (p_i, p_{i+1}] \\ 1, & \text{else.} \end{cases}$$

Then from Eq. (6) and Eq. (10), the difference  $\mathbf{c} \cdot \mathbf{m}^{(e)} - \mathbf{c}' \cdot \mathbf{m}^{(e)}$  is given by

$$\mathbf{c} \cdot \mathbf{m}^{(e)} - \mathbf{c}' \cdot \mathbf{m}^{(e)} = \sum_{j=0}^{6k} \left( \sum_{i=p_j+1}^{p_{j+1}} |s_i| i^e \right) x_j \quad (12)$$

Let  $A$  be a  $6k+1 \times 6k+1$  matrix with entries defined by  $A_{e,j} = \sum_{i=p_{j-1}+1}^{p_j} |s_i| i^{e-1}$  for  $e, j \in [1, 6k+1]$ . If  $\mathbf{c} \cdot \mathbf{m}^{(e)} = \mathbf{c}' \cdot \mathbf{m}^{(e)}$  for  $e \in [0, 6k]$ , we have the following linear equation

$$\begin{aligned} A\mathbf{x} &= \begin{bmatrix} \sum_{i=p_0+1}^{p_1} |s_i| i^0 & \dots & \sum_{i=p_{6k}+1}^{p_{6k+1}} |s_i| i^0 \\ \vdots & \ddots & \vdots \\ \sum_{i=p_0+1}^{p_1} |s_i| i^{6k} & \dots & \sum_{i=p_{6k}+1}^{p_{6k+1}} |s_i| i^{6k} \end{bmatrix} \begin{bmatrix} x_0 \\ \vdots \\ x_{6k} \end{bmatrix} \\ &= 0, \end{aligned} \quad (13)$$

with a solution  $x_i \in \{-1, 1\}$  for  $i \in [0, 6k]$ . We show that this is impossible unless  $A$  is a zero matrix. Suppose on the contrary that  $A$  is nonzero, let  $j_1 < \dots < j_Q$  be the indices of all nonzero columns of  $A$ . Let  $A^*$  be a submatrix of  $A$ , obtained by choosing the intersection of the first  $Q$  rows and columns with indices  $j_1, \dots, j_Q$ . Then taking the first  $Q$  linear equations from the equation set (14) and noting that the nonzero columns in  $A$  are the  $j_1, \dots, j_Q$ -th columns, we have that

$$\begin{aligned} &A^* \mathbf{x}' \\ &= \begin{bmatrix} \sum_{i=p_{j_1-1}+1}^{p_{j_1}} |s_i| i^0 & \dots & \sum_{i=p_{j_Q-1}+1}^{p_{j_Q}} |s_i| i^0 \\ \vdots & \ddots & \vdots \\ \sum_{i=p_{j_1-1}+1}^{p_{j_1}} |s_i| i^{Q-1} & \dots & \sum_{i=p_{j_Q-1}+1}^{p_{j_Q}} |s_i| i^{Q-1} \end{bmatrix} \begin{bmatrix} x_{j_1} \\ \vdots \\ x_{j_Q} \end{bmatrix} \\ &= 0 \end{aligned} \quad (14)$$

The determinant of  $A^*$  is given by

$$\begin{aligned} &\det(A^*) \\ &= \det \begin{pmatrix} \sum_{i=p_{j_1-1}+1}^{p_{j_1}} |s_i| i^0 & \dots & \sum_{i=p_{j_Q-1}+1}^{p_{j_Q}} |s_i| i^0 \\ \vdots & \ddots & \vdots \\ \sum_{i=p_{j_1-1}+1}^{p_{j_1}} |s_i| i^{Q-1} & \dots & \sum_{i=p_{j_Q-1}+1}^{p_{j_Q}} |s_i| i^{Q-1} \end{pmatrix} \\ &\stackrel{(a)}{=} \sum_{\substack{i_1 \in (p_{j_1-1}, p_{j_1}], \dots, \\ i_Q \in (p_{j_Q-1}, p_{j_Q}]}} \det \begin{pmatrix} |s_{i_1}| i_1^0 & \dots & |s_{i_Q}| i_Q^0 \\ \vdots & \ddots & \vdots \\ |s_{i_1}| i_1^{Q-1} & \dots & |s_{i_Q}| i_Q^{Q-1} \end{pmatrix} \\ &\stackrel{(b)}{=} \sum_{\substack{i_1 \in (p_{j_1-1}, p_{j_1}], \dots, \\ i_Q \in (p_{j_Q-1}, p_{j_Q}]}} \left[ \prod_{q=1}^Q |s_{i_q}| \det \begin{pmatrix} i_1^0 & \dots & i_Q^0 \\ \vdots & \ddots & \vdots \\ i_1^{Q-1} & \dots & i_Q^{Q-1} \end{pmatrix} \right] \\ &\stackrel{(c)}{=} \sum_{\substack{i_1 \in (p_{j_1-1}, p_{j_1}], \dots, \\ i_Q \in (p_{j_Q-1}, p_{j_Q}]}} \left[ \prod_{q=1}^Q |s_{i_q}| \prod_{1 \leq m < \ell \leq Q} (i_\ell - i_m) \right], \end{aligned} \quad (15)$$

where equality (a) follows from the multi-linearity of the determinant

$$\begin{aligned} &\det([\mathbf{v}_1 \dots a\mathbf{v}_i + b\mathbf{v} \dots \mathbf{v}_q]) \\ &= a \det([\mathbf{v}_1 \dots \mathbf{v}_i \dots \mathbf{v}_q]) \\ &\quad + b \det([\mathbf{v}_1 \dots \mathbf{v}_{i-1} \mathbf{v} \mathbf{v}_{i+1} \dots \mathbf{v}_q]) \end{aligned}$$

for any integers  $q$  and  $i$  and  $q$ -dimensional vectors  $\mathbf{v}_1, \dots, \mathbf{v}_q, \mathbf{v}$ . Equality (b) follows from the linearity of the determinant

$$\det([\mathbf{v}_1 \dots a\mathbf{v}_i \dots \mathbf{v}_q]) = a \det([\mathbf{v}_1 \dots \mathbf{v}_i \dots \mathbf{v}_q])$$



for any integers  $q$  and  $i$  and  $q$ -dimensional vectors  $\mathbf{v}_1, \dots, \mathbf{v}_q$ . Equality (c) follows from the determinant of Vandermonde matrix

$$\det \begin{pmatrix} i_1^0 & \dots & i_q^0 \\ \vdots & \ddots & \vdots \\ i_1^{q-1} & \dots & i_q^{q-1} \end{pmatrix} = \prod_{1 \leq m < \ell \leq q} (i_\ell - i_m)$$

for any integers  $q, i_1, \dots, i_q$ . The determinant  $\det(A^*)$  is positive since  $i_\ell > i_m$  for  $\ell > m$ . and for  $i_1 \in (p_{j_1-1}, p_{j_1}], \dots, i_q \in (p_{j_Q-1}, p_{j_Q}]$ . Note that all the columns of  $A^*$  are nonzero. Therefore, the linear equation  $A^* \mathbf{x}' = 0$  does not have nonzero solutions, contradicting the fact that  $\mathbf{x}' = (x_{j_1}, \dots, x_{j_Q}) \in \{-1, 1\}^Q$ . Hence  $A$  is a zero matrix, meaning that

$$|S_1 \cap [i, n]| - |S_2 \cap [i, n]| + |S_1^c \cap [i, n]| - |S_2^c \cap [i, n]| \\ = |\Delta \cap [i, n]| - |\Delta' \cap [i, n]| = 0$$

for  $i \in \{1, \dots, n\}$ . This implies  $\Delta = \Delta'$  and thus  $\mathbf{c} = \mathbf{c}'$ . Hence Proposition 2 is proved.  $\square$

### A. Proof of Lemma 1

We are now ready to prove Lemma 1, which states that  $\mathbb{1}_{\text{sync}}(\mathbf{c}) = \mathbb{1}_{\text{sync}}(\mathbf{c}')$  for sequences  $\mathbf{c}$  and  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$  satisfying  $f(\mathbb{1}_{\text{sync}}(\mathbf{c})) = f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))$ . From Proposition 1 we have that  $\mathbb{1}_{\text{sync}}(\mathbf{c}') \in \mathcal{B}_{3k}(\mathbb{1}_{\text{sync}}(\mathbf{c}))$ . Then, it is not hard to see that  $(\mathbb{1}_{\text{sync}}(\mathbf{c}')_i, \dots, \mathbb{1}_{\text{sync}}(\mathbf{c}')_n) \in \mathcal{B}_{3k}((\mathbb{1}_{\text{sync}}(\mathbf{c})_i, \dots, \mathbb{1}_{\text{sync}}(\mathbf{c})_n))$ . This implies that  $||\Delta \cap [i, n]| - |\Delta' \cap [i, n]| \leq 3k$ , where  $\Delta = \{i : \mathbb{1}_{\text{sync}}(\mathbf{c})_i = 1\}$  and  $\Delta' = \{i : \mathbb{1}_{\text{sync}}(\mathbf{c}')_i = 1\}$ . According to the forth line in Eq. (6), we have that

$$\begin{aligned} & |\mathbb{1}_{\text{sync}}(\mathbf{c}) \cdot \mathbf{m}^{(e)} - \mathbb{1}_{\text{sync}}(\mathbf{c}') \cdot \mathbf{m}^{(e)}| \\ &= \left| \sum_{i=1}^n (|\Delta \cap [i, n]| - |\Delta' \cap [i, n]|) i^e \right| \\ &\leq \sum_{i=1}^n 3ki^e \\ &< 3kn^{e+1}. \end{aligned} \quad (16)$$

If  $f(\mathbb{1}_{\text{sync}}(\mathbf{c})) = f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))$ , then

$$\mathbb{1}_{\text{sync}}(\mathbf{c}) \cdot \mathbf{m}^{(e)} \equiv \mathbb{1}_{\text{sync}}(\mathbf{c}') \cdot \mathbf{m}^{(e)} \pmod{3kn^{e+1}} \quad (17)$$

for  $e \in [0, 6k]$ . Equations (17) and (16) imply that  $\mathbb{1}_{\text{sync}}(\mathbf{c}) \cdot \mathbf{m}^{(e)} = \mathbb{1}_{\text{sync}}(\mathbf{c}') \cdot \mathbf{m}^{(e)}$  for  $e \in [0, 6k]$ . Since  $\mathbb{1}_{\text{sync}}(\mathbf{c}') \in \mathcal{B}_{3k}(\mathbb{1}_{\text{sync}}(\mathbf{c}))$  and  $\mathbb{1}_{\text{sync}}(\mathbf{c}), \mathbb{1}_{\text{sync}}(\mathbf{c}') \in \mathcal{R}_{3k}$ , from Proposition 2 we conclude that  $\mathbb{1}_{\text{sync}}(\mathbf{c}) = \mathbb{1}_{\text{sync}}(\mathbf{c}')$ . Hence Lemma 1 is proved.

### B. Proof of Lemma 2

Based on Lemma 1, we now show Lemma 2. Specifically, we show that there exists a function  $p : \{0, 1\}^n \rightarrow [1, 2^{2k \log n + o(\log n)}]$  such that  $\mathbb{1}_{\text{sync}}(\mathbf{c}) = \mathbb{1}_{\text{sync}}(\mathbf{c}')$  for sequences  $\mathbf{c}$  and  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$  satisfying  $(f(\mathbb{1}_{\text{sync}}(\mathbf{c})) \bmod p(\mathbf{c}), p(\mathbf{c})) = (f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) \bmod p(\mathbf{c}'), p(\mathbf{c}'))$ .

Lemma 1 implies that  $f(\mathbb{1}_{\text{sync}}(\mathbf{c})) \neq f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))$  for  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c}) \setminus \{\mathbf{c}\}$ , if  $\mathbb{1}_{\text{sync}}(\mathbf{c}) \neq \mathbb{1}_{\text{sync}}(\mathbf{c}')$ . Hence

$|f(\mathbb{1}_{\text{sync}}(\mathbf{c})) - f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))| \neq 0$  for  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c}) \setminus \{\mathbf{c}\}$ , where  $f(\mathbb{1}_{\text{sync}}(\mathbf{c}))$  and  $f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))$  denote the integer presentation of their vector form. The integers are in the range  $[0, (3k)^{6k+1} n^{(3k+1)(6k+1)} - 1]$ . According to Lemma 7, the number of divisors of  $|f(\mathbb{1}_{\text{sync}}(\mathbf{c})) - f(\mathbb{1}_{\text{sync}}(\mathbf{c}'))|$  is upper bounded by

$$2^{2[(3k+1)(6k+1) \ln n + (6k+1) \ln 3k] / \ln((3k+1)(6k+1) \ln n + (6k+1) \ln 3k)} \\ = 2^{o(\log n)},$$

where the equality holds since  $k = o(\sqrt{\log \log n})$ . For any sequence  $\mathbf{c} \in \{0, 1\}^n$ , let

$$\mathcal{P}(\mathbf{c}) = \{p : p \text{ divides } |f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) - f(\mathbb{1}_{\text{sync}}(\mathbf{c}))| \\ \text{for some } \mathbf{c}' \in \mathcal{B}_k(\mathbf{c}) \setminus \{\mathbf{c}\} \text{ such that } \mathbb{1}_{\text{sync}}(\mathbf{c}) \neq \mathbb{1}_{\text{sync}}(\mathbf{c}')\}$$

be the set of all divisors of the numbers  $\{|f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) - f(\mathbb{1}_{\text{sync}}(\mathbf{c}))| : \mathbf{c}' \in \mathcal{B}_k(\mathbf{c}) \setminus \{\mathbf{c}\} \text{ and } \mathbb{1}_{\text{sync}}(\mathbf{c}) \neq \mathbb{1}_{\text{sync}}(\mathbf{c}')\}$ . Since  $|\mathcal{B}_k(\mathbf{c})| \leq \binom{n}{k} 2^k \leq 2n^{2k}$ , we have that

$$|\mathcal{P}(\mathbf{c})| \leq 2n^{2k} 2^{o(\log n)} \\ = 2^{2k \log n + o(\log n)}.$$

Therefore, there exists a number  $p(\mathbf{c}) \in [1, 2^{2k \log n + o(\log n)}]$  such that  $p(\mathbf{c})$  does not divide  $|f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) - f(\mathbb{1}_{\text{sync}}(\mathbf{c}))|$  for all  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c}) \setminus \{\mathbf{c}\}$  satisfying  $\mathbb{1}_{\text{sync}}(\mathbf{c}) \neq \mathbb{1}_{\text{sync}}(\mathbf{c}')$ . Hence, if  $f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) \equiv f(\mathbb{1}_{\text{sync}}(\mathbf{c})) \pmod{p(\mathbf{c})}$  and  $\mathbf{c}' \in \mathcal{B}_k(\mathbf{c})$ , we have that  $p(\mathbf{c})$  divides  $|f(\mathbb{1}_{\text{sync}}(\mathbf{c}')) - f(\mathbb{1}_{\text{sync}}(\mathbf{c}))|$ , and thus that  $\mathbb{1}_{\text{sync}}(\mathbf{c}') = \mathbb{1}_{\text{sync}}(\mathbf{c})$ . This completes the proof of Lemma 2.

## V. HASH FOR $k$ -DENSE SEQUENCES

In this section, we present a hash function of size  $4k \log n + o(\log n)$  bits for correcting  $k$  deletions in a  $k$ -dense sequence  $\mathbf{c}$ , when the synchronization vector  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  is known. This proves Lemma 3. Recall that a sequence  $\mathbf{c}$  is  $k$ -dense if there is a 0 run of length at most  $L$  between any two 1's in the synchronization vector  $\mathbb{1}_{\text{sync}}(\mathbf{c})$ , where  $L$  is given in (5).

Let the indices of the 1 entries in  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  be  $t_1 < t_2 < \dots < t_J$ , where  $J = \sum_{i=1}^n \mathbb{1}_{\text{sync}}(\mathbf{c})_i$  is the number of 1 entries in  $\mathbb{1}_{\text{sync}}(\mathbf{c})$ . For notation convenience, let  $t_0 = 0$  and  $t_{J+1} = n + 1$ . Split  $\mathbf{c}$  into blocks  $\mathbf{a}_0, \dots, \mathbf{a}_J$ , where

$$\mathbf{a}_j = (c_{t_j+1}, c_{t_j+2}, \dots, c_{t_{j+1}-1}) \quad (18)$$

for  $j \in [0, J]$ . The blocks are separated by the synchronization patterns. Since  $\mathbf{c}$  is  $k$ -dense, the length  $|\mathbf{a}_j|$  of  $\mathbf{a}_j$  is at most  $L$ . The goal is to protect all blocks  $\mathbf{a}_j$ ,  $j \in [0, J]$  and then  $\mathbf{c}$ .

We will show that given  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  and a length  $n - k$  subsequence  $\mathbf{d}$  of  $\mathbf{c}$ , most of the blocks  $\mathbf{a}_j$  can be recovered with up to  $2k$  block errors. This is done by noticing that no deletion occurs in most of the blocks and their boundary, which are marked by the synchronization patterns in  $\mathbf{c}$ . These blocks with no deletions inside are not destructed and appear in  $\mathbf{d}$  with bits indices decrease by an integer at most  $k$ . They can be identified by looking at the synchronization patterns in  $\mathbf{d}$ .

For blocks that are not recovered, we show that they can be recovered with up to  $k$  deletion errors. Then we use the  $k$ -deletion correcting hash function  $H(\mathbf{a}_j)$ ,  $j \in [0, J]$  defined



in at most  $2k$  block errors  $\mathbf{a}'_j \neq \mathbf{a}_j$ . Therefore, the sequence  $\mathbf{c}$  can be recovered.

## VI. GENERATING $k$ -DENSE SEQUENCES

In this section we present an algorithm to compute the map  $T(\mathbf{c})$ , which transforms any sequence  $\mathbf{c} \in \{0,1\}^n$  into a  $k$ -dense sequence, and thus proves Lemma 4. Let  $\mathbf{1}^x$  and  $\mathbf{0}^y$  denote sequences of consecutive  $x$  1's and consecutive  $y$  0's, respectively. We first show in Proposition 3 that any sequence  $\mathbf{c}$  satisfying the following two properties is a  $k$ -dense sequence. Then, the algorithm for computing  $T(\mathbf{c})$  can be decomposed into two parts. In the first part, we generate a sequence  $T_1(\mathbf{c})$  that satisfies Property 1. In the second part, we use  $T_1(\mathbf{c})$  to compute  $T(\mathbf{c})$  that satisfies both properties.

**Property 1.** Every length  $B \triangleq (\lceil \log k \rceil + 5)2^{\lceil \log k \rceil + 9} \lceil \log n \rceil$  interval of  $\mathbf{c}$  contains the pattern  $\mathbf{1}^{\lceil \log k \rceil + 5}$ , i.e., for any integer  $i \in [1, n - B + 1]$ , there exists an integer  $j \in [i, i + B - \lceil \log k \rceil - 5]$  such that  $(c_j, c_{j+1}, \dots, c_{j+\lceil \log k \rceil + 4}) = \mathbf{1}^{\lceil \log k \rceil + 5}$ .

**Property 2.** Every length  $R \triangleq (3k + \lceil \log k \rceil + 4)(\lceil \log n \rceil + 9 + \lceil \log k \rceil)$  interval of  $\mathbf{c}$  contains a length  $3k + \lceil \log k \rceil + 4$  subinterval that does not contain the pattern  $\mathbf{1}^{\lceil \log k \rceil + 5}$ , i.e., for any integer  $i \in [1, n - R + 1]$ , there exists an integer  $j \in [i, i + R - 3k - \lceil \log k \rceil - 4]$ , such that  $(c_m, c_{m+1}, \dots, c_{m+\lceil \log k \rceil + 4}) \neq \mathbf{1}^{\lceil \log k \rceil + 5}$  for every  $m \in [j, j + 3k - 1]$ .

**Proposition 3.** If a sequence  $\mathbf{c}$  satisfies Property 1 and Property 2, then it is a  $k$ -dense sequence.

*Proof.* Let the locations of the 1 entries in  $\mathbb{1}_{\text{sync}}(\mathbf{c})$  be  $t_1 < \dots < t_J$ . Let  $t_0 = 0$  and  $t_{J+1} = n + 1$ . From Definition 2, it suffices to show that  $t_{i+1} - t_i \leq B + R + 1 = L + 1$  for any  $i \in [0, J]$ .

According to Property 2, there exists an index  $j^* \in [t_i, t_i + R - 3k - \lceil \log k \rceil - 4]$ , such that  $(c_m, c_{m+1}, \dots, c_{m+\lceil \log k \rceil + 4}) \neq \mathbf{1}^{\lceil \log k \rceil + 5}$  for every  $m \in [j^*, j^* + 3k - 1]$ . According to Property 1, there exists an integer  $x \in [j^* + 1, j^* + B]$  such that  $(c_x, c_{x+1}, \dots, c_{x+\lceil \log k \rceil + 4}) = \mathbf{1}^{\lceil \log k \rceil + 5}$ . Let  $\ell = \min\{x \geq j^* : (c_x, c_{x+1}, \dots, c_{x+\lceil \log k \rceil + 4}) = \mathbf{1}^{\lceil \log k \rceil + 5}\}$ . Then we have that  $\ell \neq j^*$ ,  $\ell \leq x \leq j^* + B$ , and thus that  $\ell \in [j^* + 1, j^* + B]$ . In addition,  $(c_m, c_{m+1}, \dots, c_{m+\lceil \log k \rceil + 4}) \neq \mathbf{1}^{\lceil \log k \rceil + 5}$  for every  $m \in [j^*, \ell) = \{j^*, \dots, \ell - 1\}$ . By definition of  $j^*$ , we have that  $\ell - j^* \geq 3k$ . Since  $(c_\ell, c_{\ell+1}, \dots, c_{\ell+\lceil \log k \rceil + 4}) = \mathbf{1}^{\lceil \log k \rceil + 5}$ , we have that  $\mathbb{1}_{\text{sync}}(\mathbf{c})_\ell = 1$ . Therefore, we conclude that

$$\begin{aligned} t_{i+1} - t_i &\leq \ell - t_i \\ &\leq j^* + B - t_i \\ &\leq R + B + 1 \\ &= L + 1 \end{aligned}$$

□

## A. Generating sequences satisfying Property 1

Given a sequence  $\mathbf{c} \in \{0,1\}^n$ , we now generate  $T_1(\mathbf{c}) \in \{0,1\}^{n+2^{\lceil \log k \rceil + 10}}$  that satisfies Property 1. The idea is to repeatedly delete the length  $B$  subsequences of  $\mathbf{c}$  that do not contain the pattern  $\mathbf{1}^{\lceil \log k \rceil + 5}$ , and append length  $B$  subsequences containing  $\mathbf{1}^{\lceil \log k \rceil + 5}$  to the end, without losing the information of the deleted subsequences. The deleting and appending procedure repeats until no length  $B$  subsequence with no  $\mathbf{1}^{\lceil \log k \rceil + 5}$  pattern is found. Notice that any binary sequence containing no pattern  $\mathbf{1}^{\lceil \log k \rceil + 5}$  can be regarded as a sequence of symbols with alphabet size  $2^{\lceil \log k \rceil + 5} - 1$ . Hence, such binary sequence can be compressed to a shorter binary sequence. In this way, we can add the index of the deleted subsequence and the pattern  $\mathbf{1}^{\lceil \log k \rceil + 5}$  to the compressed sequence. The sequence keeps the same size after the deleting and appending procedure.

Note that the above procedure keeps appending length  $B$  subsequences to the end. Hence the suffix of  $T_1(\mathbf{c})$  are appended bits. To guarantee that these appended bits are not deleted in the procedure, we keep track of the end index of the non-appended bits  $n'$  and always delete the bits with indices at most  $n'$ . To deal with cases when a length  $B$  subsequence to be deleted overlaps with the appended bits, we delete only non-appended bits from it. Then, we append shorter subsequences such that the total length does not change. The decoder detects the shorter appended subsequences by looking at the length of the 1 run in it.

Before presenting the details of encoding and decoding, we need the following proposition, which states that the length  $B$  binary sequences containing no  $\mathbf{1}^{\lceil \log k \rceil + 5}$  pattern can be compressed to shorter binary sequences.

**Proposition 4.** Let  $S$  be the set of sequences  $\mathbf{b} \in \{0,1\}^B$  such that  $(b_i, \dots, b_{i+\lceil \log k \rceil + 4}) \neq \mathbf{1}^{\lceil \log k \rceil + 5}$  for every  $i \in [1, B - \lceil \log k \rceil - 4]$ . There exists an invertible map  $\phi : S \rightarrow \{0,1\}^{B - \lceil \log n \rceil - 2^{\lceil \log k \rceil + 12}}$ , such that both  $\phi$  and its inverse  $\phi^{-1}$  can be computed in  $O(B)$  time.

*Proof.* For any  $\mathbf{b} \in S$ , the procedure for computing  $\phi(\mathbf{b})$  is as follows. Split  $\mathbf{b}$  into  $2^{\lceil \log k \rceil + 9} \lceil \log n \rceil$  blocks of length  $(\lceil \log k \rceil + 5)$ . Since each block is not  $\mathbf{1}^{\lceil \log k \rceil + 5}$ , it can be represented by a symbol of alphabet size  $2^{\lceil \log k \rceil + 5} - 1$ . Therefore, the sequence  $\mathbf{b}$  can be uniquely represented by a sequence  $\mathbf{v}$  of  $2^{\lceil \log k \rceil + 9} \lceil \log n \rceil$  symbols, each having alphabet size  $2^{\lceil \log k \rceil + 5} - 1$ . Convert  $\mathbf{v}$  into a binary sequence  $\phi(\mathbf{b})$ . Then  $\phi(\mathbf{b})$  can be represented by a binary sequence with length

$$\begin{aligned} &\lceil \log_2[(2^{\lceil \log k \rceil + 5} - 1)2^{\lceil \log k \rceil + 9} \lceil \log n \rceil] \rceil \\ &= \lceil \log_2[(1 - 1/2^{\lceil \log k \rceil + 5})2^{\lceil \log k \rceil + 9} \lceil \log n \rceil] \rceil \\ &\quad + (\lceil \log k \rceil + 5)2^{\lceil \log k \rceil + 9} \lceil \log n \rceil \\ &\leq 16 \lceil \log n \rceil \log_2[(1 - 1/2^{\lceil \log k \rceil + 5})2^{\lceil \log k \rceil + 5}] + B + 1 \\ &\stackrel{(a)}{\leq} -16 \lceil \log n \rceil \log_2 e + B + 1 \\ &\leq B - 16 \lceil \log n \rceil + 1 \end{aligned}$$

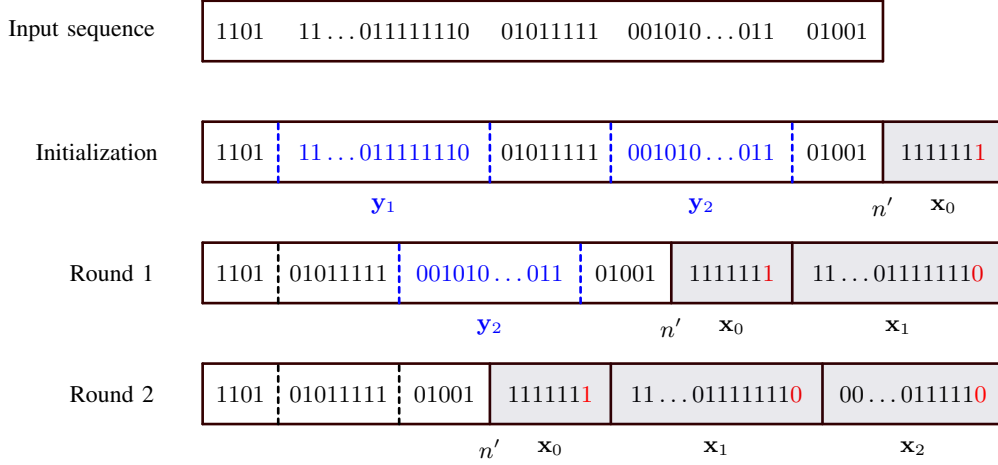


Fig. 3. This figure shows an example of how the encoding in Proposition 5 proceeds. The subsequences  $\mathbf{x}_0$ ,  $\mathbf{x}_1$ , and  $\mathbf{x}_2$  are the subsequences appended in the Initialization step, the first, and the second round of delete and insert procedure, respectively. The subsequence  $\mathbf{y}_1$  and  $\mathbf{y}_2$  are the subsequences deleted in the first and the second round of delete and insert procedure, respectively. The integer  $n'$  indicates the end point of the non-appended bits.

$$\leq B - \lceil \log n \rceil - 2\lceil \log k \rceil - 12,$$

where (a) follows from the fact that the function  $(1 - 1/x)^x$  is increasing in  $x$  for  $x > 1$  and that  $\lim_{x \rightarrow \infty} (1 - 1/x)^x = 1/e$ . Therefore,  $\phi(\mathbf{b})$  can be represented by  $B - \lceil \log n \rceil - 2\lceil \log k \rceil - 12$  bits. The inverse map  $\phi^{-1}$  can be computed by converting  $\phi(\mathbf{b})$  back to a length  $2^{\lceil \log k \rceil + 9} \lceil \log n \rceil$  sequence  $\mathbf{v}$  of alphabet size  $2^{\lceil \log k \rceil + 5} - 1$ . Then, concatenate the binary representation of symbols in  $\mathbf{v}$ , we obtain  $\mathbf{b}$ .

The complexity for computing  $\phi$  or  $\phi^{-1}$  is that of converting binary sequences to sequences of alphabet size  $2^{\lceil \log k \rceil + 5} - 1$  or vice versa, which is  $O(B)$ .  $\square$

With the function  $\phi$  defined in Proposition 4, the following proposition provides details of the encoding/decoding for computing  $T_1(\mathbf{c})$ . An example illustrating the encoding in Proposition 5 is presented in Fig. 3.

**Proposition 5.** For integers  $k$  and  $n > k$ , there exists an invertible map  $T_1 : \{0, 1\}^n \rightarrow \{0, 1\}^{n+2\lceil \log k \rceil + 10}$ , computable in  $O(n^2 k \log n \log^2 k)$  time, such that  $T_1(\mathbf{c})$  satisfies Property 1. Moreover, either

$$(T_1(\mathbf{c})_{n+\lceil \log k \rceil + 6}, \dots, T_1(\mathbf{c})_{n+2\lceil \log k \rceil + 10}) = \mathbf{1}^{\lceil \log k \rceil + 5} \text{ or } (T_1(\mathbf{c})_{n+\lceil \log k \rceil + 5}, \dots, T_1(\mathbf{c})_{n+2\lceil \log k \rceil + 9}) = \mathbf{1}^{\lceil \log k \rceil + 5}.$$

*Proof.* For a sequence  $\mathbf{c} \in \{0, 1\}^n$ , the encoding procedure for computing  $T_1(\mathbf{c})$  is as follows.

- 1) **Initialization:** Let  $T_1(\mathbf{c}) = \mathbf{c}$ . Append  $\mathbf{1}^{2\lceil \log k \rceil + 10}$  to the end of the sequence  $T_1(\mathbf{c})$ . Let  $n' = n$ . Go to Step 1.
- 2) **Step 1:** If there exists an integer  $i \in [1, n']$  such that  $(T_1(\mathbf{c})_i, T_1(\mathbf{c})_{i+1}, \dots, T_1(\mathbf{c})_{i+\lceil \log k \rceil + 4}) \neq \mathbf{1}^{\lceil \log k \rceil + 5}$  for every  $j \in [i, i + B - \lceil \log k \rceil - 5]$ , go to Step 2. Else go to Step 4.
- 3) **Step 2:** If  $i > n' - B + 1$ , go to Step 3. Else, delete  $(T_1(\mathbf{c})_i, \dots, T_1(\mathbf{c})_{i+B-1})$  from  $T_1(\mathbf{c})$  and append  $(i, \phi(T_1(\mathbf{c})_i, \dots, T_1(\mathbf{c})_{i+B-1}), 0, \mathbf{1}^{2\lceil \log k \rceil + 10}, 0)$  to the end of  $T_1(\mathbf{c})$ , where the appended  $i$  is encoded by  $\lceil \log n \rceil$  binary bits. Let  $n' = n' - B$ . Go to Step 1.

- 4) **Step 3:** Delete  $(T_1(\mathbf{c})_i, \dots, T_1(\mathbf{c})_{n'})$  from  $T_1(\mathbf{c})$  and append  $(i, \phi(T_1(\mathbf{c})_i, \dots, T_1(\mathbf{c})_{n'}), \mathbf{0}^{i+B-n'-1}, 0, \mathbf{1}^{2\lceil \log k \rceil + 10 - (i+B-n'-1)}, 0)$  to the end of  $T_1(\mathbf{c})$ . Let  $n' = i - 1$ , where the appended  $i$  is encoded by  $\lceil \log n \rceil$  binary bits. Go to Step 1.
- 5) **Step 4:** Output  $T_1(\mathbf{c})$ .

In the encoding procedure, Step 2 and Step 3 are the deleting and appending operation described above, while Step 3 deals with the case when the deleted subsequence overlaps with the appended bits. Note that the index of the deleted subsequence is provided in the appended subsequence, so that the decoder can recover the deleted subsequence, given the appended subsequence. According to Proposition 4 and the fact that the index  $i$  has size  $\lceil \log n \rceil$  bits, the lengths of the deleted and appended subsequences in Step 2 or Step 3 are equal. Hence, the sequence  $T_1(\mathbf{c})$  keeps constant and is  $n + 2\lceil \log k \rceil + 10$ .

We first show that the integer  $n'$  is the split index of appended bits and non-appended bits. Specifically,  $(T_1(\mathbf{c})_{n'+1}, \dots, T_1(\mathbf{c})_{n+2\lceil \log k \rceil + 10})$  are the appended bits and  $(T_1(\mathbf{c})_1, \dots, T_1(\mathbf{c})_{n'})$  are non-appended bits that remain after deleting operations in Step 2 and Step 3. In addition,  $(T_1(\mathbf{c})_{n'+1}, \dots, T_1(\mathbf{c})_{n'+2\lceil \log k \rceil + 10}) = \mathbf{1}^{2\lceil \log k \rceil + 10}$  are the bits appended in the Initialization step. The claim holds in the Initialization step. Note that in each round of Step 2 or Step 3, the deleted bits have indices at most  $n'$  and the integer  $n'$  decreases by the length of deleted subsequence. Hence, the claim always holds.

We now show that the output sequence  $T_1(\mathbf{c})$  satisfies Property 1.

We have shown that  $(T_1(\mathbf{c})_{n'+1}, \dots, T_1(\mathbf{c})_{n'+\lceil \log k \rceil + 5}) = \mathbf{1}^{\lceil \log k \rceil + 5}$ . Then according to the if conditions in Step 1 and Step 2 that lead to Step 3, the integer  $i$  in Step 3 satisfies  $1 \leq i + B - n' - 1 \leq \lceil \log k \rceil + 4$ . Otherwise the subsequence  $(T_1(\mathbf{c})_i, \dots, T_1(\mathbf{c})_{i+B-1})$  contains  $(T_1(\mathbf{c})_{n'+1}, \dots, T_1(\mathbf{c})_{n'+\lceil \log k \rceil + 5}) = \mathbf{1}^{\lceil \log k \rceil + 5}$ , which does not satisfy the if condition in Step 1. Therefore, the 1



run  $1^{2\lceil \log k \rceil + 10 - (i+B-n'-1)}$  in the appended subsequence in Step 3 has length at least  $\lceil \log k \rceil + 6$ . Thus, the  $1^{\lceil \log k \rceil + 5}$  pattern appears in the subsequence appended in each round of Step 2 or Step 3. Moreover, the index distance between the two  $1^{\lceil \log k \rceil + 5}$  patterns in two consecutively appended subsequences is at most  $B - \lceil \log k \rceil - 5$ . We conclude that for  $i' > n'$ , any length  $B$  subsequence  $(T_1(\mathbf{c})_i, \dots, T_1(\mathbf{c})_{i+B-1})$  contains the  $1^{\lceil \log k \rceil + 5}$  pattern. Furthermore, note that for any  $i \in [1, n']$ , there exists some  $j \in [i, i+B - \lceil \log k \rceil - 5]$  such that  $T_1(\mathbf{c})_j = T_1(\mathbf{c})_{j+1} = \dots = T_1(\mathbf{c})_{j+\lceil \log k \rceil + 4} = 1$ . Otherwise  $T_1(\mathbf{c})_i$  is deleted in Step 2 or Step 3. Hence, the sequence  $T_1(\mathbf{c})$  satisfies Property 1.

Note that the integer  $n'$  decreases in each round. Hence, the algorithm terminates within  $O(n)$  rounds of Step 1, Step 2 and Step 3. Therefore, the search for the  $i$  satisfying the if condition in Step 1 takes  $O(nB \log k)$  time. In addition, the deleting and appending operation in Step 2 or Step 3 takes at most  $O(n+B)$  time. Hence, the total complexity is  $O(n^2 k \log n \log^2 k)$ .

Next, we show that either  $(T_1(\mathbf{c})_{n+\lceil \log k \rceil + 6}, \dots, T_1(\mathbf{c})_{n+2\lceil \log k \rceil + 10}) = 1^{\lceil \log k \rceil + 5}$  or  $(T_1(\mathbf{c})_{n+\lceil \log k \rceil + 5}, \dots, T_1(\mathbf{c})_{n+2\lceil \log k \rceil + 9}) = 1^{\lceil \log k \rceil + 5}$ . The former holds when the appending operation only occurs in the Initialization step. Note that all subsequence appended in Step 2 or Step 3 ends with a 1 run with length at least  $\lceil \log k \rceil + 6$  followed by a 0 bit. Hence the latter holds when the appending operation in Step 2 or Step 3 occurs.

The decoding follows a reverse procedure of the encoding, by repeatedly removing the appended subsequences and use them to recover the deleted subsequences. Since the appended subsequences contain the  $\phi$  function and position of the deleted subsequences, the deleted subsequences can be recovered. The decoding stops when the end of the appended sequence becomes a 1 bit. The decoder determines the length of the appended subsequence by looking at the 1 run before the ending 0 bit. The decoding procedure that recovers  $\mathbf{c}$  from  $T_1(\mathbf{c})$  is given as follows.

- 1) **Initialization:** Let  $\mathbf{c} = T_1(\mathbf{c})$  and go to Step 1.
- 2) **Step 1:** If  $c_{n+2\lceil \log k \rceil + 10} = 0$ , find the length  $\ell$  of the 1 run that ends with  $c_{n+2\lceil \log k \rceil + 9}$ . Let  $i$  be the integer representation of  $(c_{n+4\lceil \log k \rceil + 21 - B - \ell}, c_{n+4\lceil \log k \rceil + 22 - B - \ell}, \dots, c_{n+4\lceil \log k \rceil + 20 - B - \ell + \lceil \log n \rceil})$ . Let  $\mathbf{b}$  be the sequence obtained by computing  $\phi^{-1}(c_{n+4\lceil \log k \rceil + 21 - B - \ell + \lceil \log n \rceil}, c_{n+4\lceil \log k \rceil + 22 - B - \ell + \lceil \log n \rceil}, \dots, c_{n+2\lceil \log k \rceil + 8 - \ell})$ , where the function  $\phi$  is defined in Proposition 4. Delete  $(c_{n+4\lceil \log k \rceil + 21 - B - \ell}, c_{n+4\lceil \log k \rceil + 22 - B - \ell}, \dots, c_{n+2\lceil \log k \rceil + 10})$  from  $\mathbf{c}$  and insert  $(\mathbf{b}_1, \dots, \mathbf{b}_{B-2\lceil \log k \rceil - 10 + \ell})$  at location  $i$  of  $\mathbf{c}$ . Repeat. Else delete  $c_{n+1}, \dots, c_{n+2\lceil \log k \rceil + 10}$  and go to Step 2.
- 3) **Step 2:** Output  $\mathbf{c}$ .

In each round of Step 1, the decoder processes an appended subsequence of length  $B$  or less. It can be seen that the appended subsequences in the encoding procedure end with a 0 bit except for the one appended in the Initialization step. Hence if the end of an appended sequence is a 1 bit, the subsequence is the  $1^{2\lceil \log k \rceil + 10}$  appended in the Initialization

step of the encoding procedure. Moreover, since  $1^{2\lceil \log k \rceil + 10}$  is the appended subsequence, the decoding procedure ends up in  $1^{2\lceil \log k \rceil + 10}$  after processing all other appended subsequences.

Note that the encoding procedure consists of a series of deleting and appending operations. The decoding procedure exactly reverses the series of operations in the encoding procedure. Let  $T_{1,i}(\mathbf{c})$ ,  $i \in [0, I]$  be the sequence  $T_1(\mathbf{c})$  obtained after the  $i$ -th deleting and appending operation in the encoding procedure, where  $I$  is the number of deleting and appending operations in total in the encoding procedure. We have that  $T_{1,0}(\mathbf{c}) = \mathbf{c}$  and that  $T_{1,I}(\mathbf{c})$  is the final output  $T_1(\mathbf{c})$ . Then the decoding procedure obtains  $T_{1,I-i}(\mathbf{c})$ ,  $i \in [0, I]$  after the  $i$ -th deleting and inserting operation. Hence we get the output  $T_{1,I-I}(\mathbf{c}) = \mathbf{c}$  in the decoding procedure.  $\square$

## B. Proof of Lemma 4

After having the sequence  $T_1(\mathbf{c}) \in \{0, 1\}^{n+2\lceil \log k \rceil + 10}$  satisfying Property 1, we now use  $T_1(\mathbf{c})$  to generate a sequence  $T(\mathbf{c}) \in \{0, 1\}^{n+3k+3\lceil \log k \rceil + 15}$  that satisfies both Property 1 and Property 2. According to Proposition 3,  $T(\mathbf{c})$  is a  $k$ -dense sequence and Lemma 4 is proved. The encoding of  $T(\mathbf{c})$  follows a similar repetitive deleting and appending procedure to the encoding in Proposition 5. To satisfy Property 2, we repeatedly look for length  $R$  subsequences, every length  $3k + \lceil \log k \rceil + 4$  subsequence of which contains  $1^{\lceil \log k \rceil + 5}$  patterns, delete most part of it, and then append a subsequence that contains a length  $3k + \lceil \log k \rceil + 4$  subsequence containing no  $1^{\lceil \log k \rceil + 5}$  pattern. The appended subsequence has all information about the deleted sequence, including the index and the bits, and has the same length as that of the deleted subsequence. To this end, we need to construct a map that compresses a length  $3k + \lceil \log k \rceil + 4$  sequence containing  $1^{\lceil \log k \rceil + 5}$  patterns to a shorter sequence. Moreover, the compressed sequence does not contain  $1^{\lceil \log k \rceil + 5}$  patterns. Then we can compress the deleted subsequence and add indices and markers as we did in Proposition 5. The end of an appended subsequence is marked by a 0 bit. We also keep track of the integer  $n'$ , which is the end of non-appended bits, to guarantee that the appended bits are not deleted.

Note that we do not delete the whole length  $R$  subsequences, every length  $3k + \lceil \log k \rceil + 4$  subsequence of which contains  $1^{\lceil \log k \rceil + 5}$  patterns. Instead, we split the length  $R$  subsequence into blocks of length  $3k + \lceil \log k \rceil + 4$ , each containing  $1^{\lceil \log k \rceil + 5}$  patterns. Then the first and the last blocks remain and the blocks in the middle are deleted. This is to keep the sequence  $T(\mathbf{c})$  satisfying Property 1 and protect the appended bits from being deleted.

The following proposition presents the compression map described above, which encodes a sequence containing  $1^{\lceil \log k \rceil + 5}$  patterns into a shorter sequence without the  $1^{\lceil \log k \rceil + 5}$  pattern. Similar to the encoding procedures that compute  $T(\mathbf{c})$  and  $T_1(\mathbf{c})$ , the algorithm for computing the compression map follows a delete and append process, which repeatedly deletes  $1^{\lceil \log k \rceil + 5}$  patterns and appends their indices and 0 runs to the end.

**Proposition 6.** For an integer  $k$ , let  $\mathbf{c} \in \{0, 1\}^{3k + \lceil \log k \rceil + 4}$  be a sequence such that  $c_i = c_{i+1} = \dots = c_{i + \lceil \log k \rceil + 4} = 1$  for some  $i \in [1, 3k]$ . There exists an invertible mapping  $T_2 : \{0, 1\}^{3k + \lceil \log k \rceil + 4} \rightarrow \{0, 1\}^{3k + \lceil \log k \rceil + 3}$ , such that  $T_2(\mathbf{c})$  contains no  $\lceil \log k \rceil + 5$  consecutive 1 bits. Both  $T_2$  and its inverse  $T_2^{-1}$  are computable in  $O(k^2 \log k)$  time.

*Proof.* Given  $\mathbf{c} \in \{0, 1\}^{3k + \lceil \log k \rceil + 4}$ , the encoding procedure for computing  $T_2(\mathbf{c})$  is as follows.

- 1) **Initialization:** Let  $T_2(\mathbf{c}) = \mathbf{c}$ . Append 0 to the end of the sequence  $T_2(\mathbf{c})$ . Find the smallest  $i \in [1, 3k]$  such that  $T_2(\mathbf{c})_i = T_2(\mathbf{c})_{i+1} = \dots = T_2(\mathbf{c})_{i + \lceil \log k \rceil + 4} = 1$ . Delete  $(T_2(\mathbf{c})_i, \dots, T_2(\mathbf{c})_{i + \lceil \log k \rceil + 4})$  from  $T_2(\mathbf{c})$  and append  $(i, \mathbf{0}^{\lceil \log k \rceil + 3 - \lceil \log(3k) \rceil})$  to the end of  $T_2(\mathbf{c})$ , where the appended  $i$  is encoded by  $\lceil \log 3k \rceil$  binary bits. Let  $n' = 3k - 1$  and  $i = 1$ . Go to Step 1.
- 2) **Step 1:** If there exists an integer  $i \leq n'$  such that  $T_2(\mathbf{c})_i = T_2(\mathbf{c})_{i+1} = \dots = T_2(\mathbf{c})_{i + \lceil \log k \rceil + 4} = 1$ , delete  $(T_2(\mathbf{c})_i, \dots, T_2(\mathbf{c})_{i + \lceil \log k \rceil + 4})$  from  $T_2(\mathbf{c})$  and append  $(i, \mathbf{0}^{\lceil \log k \rceil + 4 - \lceil \log(3k) \rceil}, 1)$  to the end of  $T_2(\mathbf{c})$ . Let  $n' = n' - \lceil \log k \rceil - 5$  and  $i = 1$ . Repeat. Else go to Step 2.
- 3) **Step 2:** Output  $T_2(\mathbf{c})$ .

There are deleting and appending operations in both the Initialization step and Step 1. In the Initialization step, the length of the deleted subsequence is larger than the length of the appended subsequence by 2. Hence after appending the 0 bit in the beginning, the length of  $T_2(\mathbf{c})$  decreases by 1 after the Initialization step. The lengths of the deleted and appended subsequence in Step 1 are equal. Hence, the length of the sequence  $T_2(\mathbf{c})$  keeps constant after the Initialization step and is

$$3k + \lceil \log k \rceil + 4 + 1 - \lceil \log k \rceil - 5 + \lceil \log k \rceil + 3 = 3k + \lceil \log k \rceil + 3.$$

We show that  $(T_2(\mathbf{c})_{n'+1}, \dots, T_2(\mathbf{c})_{3k + \lceil \log k \rceil + 3})$  are appended bits and  $(T_2(\mathbf{c})_1, \dots, T_2(\mathbf{c})_{n'})$  are non-appended bits remain in  $T_2(\mathbf{c})$  after deletions. In particular,  $T_2(\mathbf{c})_{n'+1} = 0$  is the bit appended in the beginning of the Initialization step. The claims hold after the Initialization step. Note that when  $T_2(\mathbf{c})_{n'} = 0$ , the integer  $i$  satisfying the if condition in Step 1 is at most  $n' - \lceil \log k \rceil - 4$ . Otherwise  $(T_2(\mathbf{c})_i, \dots, T_2(\mathbf{c})_{i + \lceil \log k \rceil + 4}) \neq \mathbf{1}$ . Therefore, the deleted bits in Step 1 have indices at most  $n'$ . Moreover, the integer  $n'$  decreases by the length of the deleted bits. Therefore, we conclude that  $(T_2(\mathbf{c})_{n'+1}, \dots, T_2(\mathbf{c})_{3k + \lceil \log k \rceil + 3})$  consists of appended bits and that  $T_2(\mathbf{c})_{n'+1} = 0$  is the bit appended at the beginning of the Initialization step. The appended bits  $(T_2(\mathbf{c})_{n'+1}, \dots, T_2(\mathbf{c})_{3k + \lceil \log k \rceil + 3})$  are not deleted in the procedure.

We now show that  $T_2(\mathbf{c})$  contains no  $\mathbf{1}^{\lceil \log k \rceil + 5}$  patterns. Note that the  $\mathbf{1}^{\lceil \log k \rceil + 5}$  patterns with indices at most  $n'$  are deleted in the encoding procedure. Since  $T_2(\mathbf{c})_{n'+1} = 0$ , there is no  $\mathbf{1}^{\lceil \log k \rceil + 5}$  pattern containing the bit  $T_2(\mathbf{c})_{n'+1}$ . Hence a  $\mathbf{1}^{\lceil \log k \rceil + 5}$  pattern has indices at least  $n' + 1$ . Moreover, since the  $\lceil \log k \rceil + 4$ -th bit in each appended subsequence is a 0 bit, there is no  $\mathbf{1}^{\lceil \log k \rceil + 5}$  pattern in the appended bits. Hence,

no  $\mathbf{1}^{\lceil \log k \rceil + 5}$  with indices at least  $n' + 1$  exists. This implies that  $T_2(\mathbf{c})$  does not contain  $\mathbf{1}^{\lceil \log k \rceil + 5}$  patterns.

Since  $n'$  decreases in each step, the algorithm terminates within  $O(k)$  iterations of Step 1. Since it takes  $O(k \log k)$  to look for the integer  $i$ , the total complexity is  $O(k^2 \log k)$ .

The decoding  $T_2^{-1}(\mathbf{c})$ , which recovers  $\mathbf{c}$  from  $T_2(\mathbf{c})$ , follows a reverse procedure of the encoding and is presented in the following.

- 1) **Initialization:** Let  $\mathbf{c} = T_2(\mathbf{c})$  and go to Step 1.
- 2) **Step 1:** If  $c_{3k + \lceil \log k \rceil + 3} = 1$ , let  $i$  be the integer representation of  $(c_{3k-1}, c_{3k}, \dots, c_{3k + \lceil \log 3k \rceil - 2})$ . Delete  $(c_{3k-1}, c_{3k}, \dots, c_{3k + \lceil \log k \rceil + 3})$  from  $\mathbf{c}$  and insert  $\mathbf{1}^{\lceil \log k \rceil + 5}$  at location  $i$  of  $\mathbf{c}$ . Repeat. Else go to Step 2.
- 3) **Step 2:** Let  $i$  be the decimal representation of  $(c_{3k+1}, c_{3k+2}, \dots, c_{3k + \lceil \log(3k) \rceil})$ . Delete  $(c_{3k}, c_{3k+2}, \dots, c_{3k + \lceil \log k \rceil + 3})$  from  $\mathbf{c}$  and insert  $\mathbf{1}^{\lceil \log k \rceil + 5}$  at location  $i$  of  $\mathbf{c}$ . Output  $\mathbf{c}$ .

Note that in the encoding procedure, the appended subsequence in the Initialization Step ends with a 0. The appended subsequence in Step 1 ends with a 1. Hence the algorithm stops when  $c_{3k + \lceil \log k \rceil + 3} = 0$  and all subsequences appended in Step 1 of the encoding have been processed.

Similarly to the proof of correctness of decoding in Proposition 5, the decoding procedure exactly reverses the series of operations in the encoding procedure. Note that the appended subsequences contain the index of the deleted  $\mathbf{1}^{\lceil \log k \rceil + 5}$  patterns. Let  $T_{2,i}(\mathbf{c})$ ,  $i \in [0, I]$  be the sequence obtained after the  $i$ -th deleting and appending operation in the encoding procedure, where  $I$  is the number of deleting and appending operations in total in the encoding procedure. Then  $T_{2,i}(\mathbf{c})$  is the sequence obtained after the  $I - i$ -th deleting and inserting operation in the decoding procedure. Therefore, the decoding procedure recovers the sequence  $\mathbf{c}$  after the  $I$ -th operation.

The complexity of the decoding has the same order  $O(k^2 \log k)$  as that of the encoding.  $\square$

We are now ready to present the encoding and decoding procedures for computing  $T(\mathbf{c})$ , which generates  $k$ -dense sequences that satisfy Property 1 and Property 2. The encoding procedure is as follows.

- 1) **Initialization:** Let  $T(\mathbf{c}) = T_1(\mathbf{c})$ . Append  $(\mathbf{0}^{3k}, \mathbf{1}^{\lceil \log k \rceil + 5})$  to the end of the sequence  $T(\mathbf{c})$ . Let  $n' = n + 2\lceil \log k \rceil + 10$  (the length of  $T_1(\mathbf{c})$ ). Go to Step 1.
- 2) **Step 1:** If there exists an integer  $i \leq \min\{n', n + 3k + 3\lceil \log k \rceil + 16 - R\}$ , such that for every  $j \in [i, i + R - 3k - \lceil \log k \rceil - 4]$ , there exists an integer  $m \in [j, j + 3k - 1]$  satisfying  $(T(\mathbf{c})_m, T(\mathbf{c})_{m+1}, \dots, T(\mathbf{c})_{m + \lceil \log k \rceil + 4}) = \mathbf{1}^{\lceil \log k \rceil + 5}$ , split  $(T(\mathbf{c})_i, T(\mathbf{c})_{i+1}, \dots, T(\mathbf{c})_{i+R-1})$  into  $(\lceil \log n \rceil + 9 + \lceil \log k \rceil)$  blocks  $\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_{\lceil \log n \rceil + 9 + \lceil \log k \rceil}$  of length  $3k + \lceil \log k \rceil + 4$ . Delete  $(\mathbf{b}_2, \dots, \mathbf{b}_{\lceil \log n \rceil + 8 + \lceil \log k \rceil})$  from  $T(\mathbf{c})$  and append  $(0, T_2(\mathbf{b}_2), T_2(\mathbf{b}_3), \dots, T_2(\mathbf{b}_{\lceil \log n \rceil + 8 + \lceil \log k \rceil}), i + 3k + \lceil \log k \rceil + 4, \mathbf{1}^{\lceil \log k \rceil + 5}, 0)$  to the end of  $T(\mathbf{c})$ , where the appended  $i + 3k + \lceil \log k \rceil + 4$  encoded by  $\lceil \log n \rceil$  binary bits. Let  $n' = n' - R + 6k + 2\lceil \log k \rceil + 8$ . Repeat. Else go to Step 2.

### 3) Step 2: Output $T(\mathbf{c})$ .

Note that the index  $i + 3k + \lceil \log k \rceil + 4$  has size  $\log n$  bits. Then from Proposition 6, it can be verified that the lengths of the deleted and appended subsequences in Step 1 are equal. Hence  $T(\mathbf{c})$  keeps constant and is  $n + 3k + 3\lceil \log k \rceil + 15$ .

Similar to the encoding in Proposition 5 and Proposition 6, we show that the integer  $n'$  marks the end of the non-appended bits, i.e., the subsequence  $(T(\mathbf{c})_{n'+1}, \dots, (\mathbf{c})_{n+3k+3\lceil \log k \rceil+15})$  consists of appended bits and  $(T(\mathbf{c})_1, \dots, T(\mathbf{c})_{n'})$  are non-appended bits that remain after deletions. In addition, we show that  $T(\mathbf{c})_{n'+1}, \dots, T(\mathbf{c})_{n'+3k} = \mathbf{0}^{3k}$  and that the appended bits are not deleted. The claims hold in the Initialization step. Suppose the claims hold in the  $r$ -th round of Step 1. Then in the  $r + 1$ -th round of Step 1, the integer  $i$  satisfying the if condition in Step 1 must be in the range  $[1, n' - R + 3k + \lceil \log k \rceil + 4]$ . Otherwise, since  $T(\mathbf{c})_{n'+1}, \dots, T(\mathbf{c})_{n'+3k} = \mathbf{0}^{3k}$  in the  $r$ -th round, we have an integer  $n' \in [i, i + R - 3k - \lceil \log k \rceil - 4]$  such that  $(T(\mathbf{c})_m, T(\mathbf{c})_{m+1}, \dots, T(\mathbf{c})_{m+\lceil \log k \rceil+4}) \neq \mathbf{1}^{\lceil \log k \rceil+5}$  for every  $m \in [n', n' + 3k - 1]$ . This contradicts to the fact that  $i$  satisfies the if condition. Moreover, note that the block  $\mathbf{b}_{\lceil \log n \rceil+9+\lceil \log k \rceil}$  in  $(T(\mathbf{c})_i, \dots, T(\mathbf{c})_{i+R-1})$  is not deleted in Step 1, which implies that the bits with indices at least  $i + R - 3k - \lceil \log k \rceil - 4 \leq n'$  are not deleted. Note that  $n'$  decreases by the length of the deleted sequence in Step 1. We conclude that  $(T(\mathbf{c})_{n'+1}, \dots, (\mathbf{c})_{n+3k+3\lceil \log k \rceil+15})$  consists of appended bits and these bits are not deleted. Specifically,  $T(\mathbf{c})_{n'+1}, \dots, T(\mathbf{c})_{n'+3k}$  are the bits appended in the Initialization step. By induction, the claims hold.

We now show by induction on the number of rounds  $r$  that  $T(\mathbf{c})$  satisfies Property 1. From Proposition 5, the initial sequence  $T(\mathbf{c}) = (T_1(\mathbf{c}), \mathbf{0}^{3k}, \mathbf{1}^{\lceil \log k \rceil+5})$  satisfies Property 1. Hence the claim holds for  $r = 0$ . Suppose after  $r$ -th round of Step 1,  $T(\mathbf{c})$  satisfies Property 1. In the  $r + 1$ -th round, the deleting operation leaves blocks  $\mathbf{b}_1$  and  $\mathbf{b}_{\lceil \log n \rceil+9+\lceil \log k \rceil}$ , which both contain  $\mathbf{1}^{\lceil \log k \rceil+5}$  as a subsequence. Hence  $T(\mathbf{c})$  satisfies Property 1 after the deletion. In addition, all appended subsequences end with a  $\mathbf{1}^{\lceil \log k \rceil+5}$  pattern or a  $\mathbf{1}^{\lceil \log k \rceil+5}$  pattern followed by a 0 bit. Note that these appended subsequences are not deleted. Hence the index distance between two  $\mathbf{1}^{\lceil \log k \rceil+5}$  patterns in the appended bits  $(T(\mathbf{c})_{n'+1}, \dots, T(\mathbf{c})_{n+3k+3\lceil \log k \rceil+15})$  is at most  $R - 6k - 2\lceil \log k \rceil - 8 \leq B - \lceil \log k \rceil - 5$ . Therefore, The sequence  $T(\mathbf{c})$  satisfies Property 1 after the appending operation.

Next, we prove that  $T(\mathbf{c})$  satisfies Property 2. According to the encoding procedure, for any  $i \in [1, \min\{n', n + 3k + 3\lceil \log k \rceil + 16 - R\}]$ , there exists some  $j \in [i, i + R - 3k - \lceil \log k \rceil - 4]$ , such that  $(T(\mathbf{c})_m, T(\mathbf{c})_{m+1}, \dots, T(\mathbf{c})_{m+\lceil \log k \rceil+4}) \neq \mathbf{1}^{\lceil \log k \rceil+5}$  for every  $m \in [j, j + 3k - 1]$ . Otherwise, the encoding does not stop. Note that the appended bits  $(T(\mathbf{c})_{n'+1}, \dots, T(\mathbf{c})_{n+3k+3\lceil \log k \rceil+15})$  are not deleted. Hence for  $i \in [n' + 1, n + 3k + 3\lceil \log k \rceil + 16 - R]$ , the interval  $[i, i + R - 1]$  contains the first  $3k + \lceil \log k \rceil + 4$  bits  $(0, T_2(\mathbf{b}_2))$  of some appended subsequence, where  $\mathbf{b}_2 \in \{0, 1\}^{3k+\lceil \log k \rceil+4}$  contains the  $\mathbf{1}^{\lceil \log k \rceil+5}$  pattern. Let  $[j, j + 3k + \lceil \log k \rceil + 3]$  be the indices of the  $3k + \lceil \log k \rceil + 4$  bits  $(0, T_2(\mathbf{b}_2))$  in  $T(\mathbf{c})$ . According to Proposition 6, the function  $T_2(\mathbf{b}_2)$  does not contain the  $\mathbf{1}^{\lceil \log k \rceil+5}$  pattern.

Hence  $(T(\mathbf{c})_m, T(\mathbf{c})_{m+1}, \dots, T(\mathbf{c})_{m+\lceil \log k \rceil+4}) \neq \mathbf{1}^{\lceil \log k \rceil+5}$  for every  $m \in [j, j + 3k - 1]$ . Therefore, any interval of length  $R$  contains a length  $3k + \lceil \log k \rceil + 4$  subsequence with no  $\mathbf{1}^{\lceil \log k \rceil+5}$  pattern. Hence, the sequence  $T(\mathbf{c})$  satisfies Property 2. According to Proposition 3, we conclude that  $T(\mathbf{c})$  satisfies Property 1 and Property 2 and is  $k$ -dense.

Since  $n'$  decreases in the encoding procedure, the procedure terminates within  $O(n)$  iterations. Each iteration takes  $O(nk \log k R)$  time to search for the integer  $i$  satisfying the if condition and  $O(\log nk^2 \log k)$  to compute the  $T_2$  functions of the blocks. Hence the complexity is at most  $O(n^2 k^2 \log k (\log n))$ . Therefore, the total complexity is  $\text{poly}(n, k)$ .

Finally we present the following decoding procedure that recovers  $\mathbf{c}$  from  $T(\mathbf{c})$ , which follows a reverse procedure of the encoding.

- 1) **Initialization:** Let  $\mathbf{c} = T(\mathbf{c})$  and go to Step 1.
- 2) **Step 1:** If  $c_{n+3k+3\lceil \log k \rceil+15} = 0$ , let  $i$  be the integer representation of  $(c_{n+3k+2\lceil \log k \rceil+10-\lceil \log n \rceil}, c_{n+3k+2\lceil \log k \rceil+11-\lceil \log n \rceil}, \dots, c_{n+3k+2\lceil \log k \rceil+9})$ . Split  $(c_{n+9k+5\lceil \log k \rceil+25-R}, \dots, c_{n+3k+2\lceil \log k \rceil+9-\lceil \log n \rceil})$  into  $\lceil \log n \rceil + 7 + \lceil \log k \rceil$  blocks  $(\mathbf{b}'_1, \dots, \mathbf{b}'_{\lceil \log n \rceil+7+\lceil \log k \rceil})$  of length  $3k + \lceil \log k \rceil + 3$ . Compute  $\mathbf{b}_j = T_2^{-1}(\mathbf{b}'_j)$  for  $j \in [1, \lceil \log n \rceil + 7 + \lceil \log k \rceil]$ , where  $T_2^{-1}(\mathbf{b}'_j)$  is defined in Proposition 6. Delete  $(c_{n+9k+5\lceil \log k \rceil+24-R}, \dots, c_{n+3k+3\lceil \log k \rceil+15})$  from  $\mathbf{c}$  and insert  $\mathbf{b}_1, \dots, \mathbf{b}_{\lceil \log n \rceil+7+\lceil \log k \rceil}$  at location  $i$  of  $\mathbf{c}$ . Repeat. Else delete  $(c_{n+2\lceil \log k \rceil+11}, \dots, c_{n+3k+3\lceil \log k \rceil+15})$  and go to Step 2.
- 3) **Step 2:** Output  $T_1^{-1}(\mathbf{c})$ .

According to the encoding procedure, the inserted bits end with a 1 entry in the Initialization Step and with a 0 entry in Step 1. Note that the inserted bits are not deleted in the encoding procedure. Hence the decoding algorithm stops when an ending 1 entry is detected.

Similar to the proof of correctness of decoding in Proposition 5 and Proposition 6, the decoding procedure exactly reverses the series of operations in the encoding procedure. Therefore, the decoding procedure decodes the sequence  $\mathbf{c}$  correctly.

## VII. CONCLUSION AND FUTURE WORK

We construct a  $k$ -deletion correcting code with optimal order redundancy. Interesting open problems include finding complexity  $O(N^{O(1)})$  encoding/decoding algorithms for our code, as well as constructing a systematic  $k$ -deletion correcting code with optimal redundancy.

## ACKNOWLEDGMENT

The authors would like to thank the anonymous reviewers for their valuable suggestions that greatly improved the clarity and the organization of the paper.

## REFERENCES

- [1] V. I. Levenshtein, "Binary codes capable of correcting deletions, insertions, and reversals," *Soviet physics doklady*, vol. 10, no. 8, pp. 707–710, 1966.

- [2] R. R. Varshamov and G. M. Tenengolts, "Codes which correct single asymmetric errors," *Autom. Remote Control*, vol. 26, no. 2, pp. 286–290, 1965.
- [3] M. Mitzenmacher, "A survey of results for deletion channels and related synchronization channels," *Probability Surveys*, vol. 6, pp. 1–33, 2009.
- [4] A. S. Helberg and H. C. Ferreira, "On multiple insertion/deletion correcting codes," *IEEE Trans. on Inf. Th.*, vol. 48, no. 1, pp. 305–308, 2002.
- [5] K. A. Abdel-Ghaffar, F. Paluncic, H. C. Ferreira, and W. A. Clarke, "On Helberg's generalization of the Levenshtein code for multiple deletion/insertion error correction," *IEEE Trans. on Inf. Th.*, vol. 58, no. 3, pp. 1804–1808, 2012.
- [6] F. Paluncic, K. A. Abdel-Ghaffar, H. C. Ferreira, and W. A. Clarke, "A multiple insertion/deletion correcting code for run-length limited sequences," *IEEE Trans. on Inf. Th.*, vol. 58, no. 3, pp. 1809–1824, 2012.
- [7] J. Brakensiek, V. Guruswami, and S. Zbarsky, "Efficient low-redundancy codes for correcting multiple deletions," *IEEE Trans. on Inf. Th.*, vol. 64, no. 5, pp. 3403–3410, 2018.
- [8] B. Haeupler, "Optimal document exchange and new codes for small number of insertions and deletions," *IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 334–347, 2019.
- [9] K. Cheng, Z. Jin, X. Li and K. Wu, "Deterministic document exchange protocols, and almost optimal binary codes for edit errors," *IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pp. 200–211, 2018.
- [10] L. J. Schulman and D. Zuckerman, "Asymptotically good codes correcting insertions, deletions, and transpositions," *IEEE Trans. on Inf. Th.*, vol. 45, no. 7, pp. 2552–2557, 1999.
- [11] V. Guruswami and C. Wang, "Deletion codes in the high-noise and high-rate regimes," *IEEE Trans. on Inf. Th.*, vol. 63, no. 4, pp. 1961–1970, 2017.
- [12] R. Gabrys and F. Sala, "Codes correcting two deletions," *IEEE Trans. on Inf. Th.*, vol. 65, no. 2, pp. 965–974, Feb. 2019.
- [13] J. Sima, N. Raviv, and J. Bruck, "Two Deletion Correcting Codes from Indicator Vectors," in *IEEE Trans. on Inf. Th.*, vol. 66, no. 4, pp. 2375–2391, 2020.
- [14] S. K. Hanna and S. El Rouayheb, "Guess & check codes for deletions, insertions, and synchronization," *IEEE Trans. on Inf. Th.*, vol. 65, no. 1, pp. 3–15, Jan. 2019.
- [15] J. L. Nicolas, "On highly composite numbers," in *Ramanujan revisited*, Urbana-Champaign, Ill., pp. 215–244, 1987.
- [16] R. Roth, *Introduction to coding theory*, Cambridge University Press, 2006.

Dr. Bruck is a recipient of the Feynman Prize for Excellence in Teaching, the Sloan Research Fellowship, the National Science Foundation Young Investigator Award, the IBM Outstanding Innovation Award and the IBM Outstanding Technical Achievement Award.

**Jin Sima** is a PhD candidate in the Department of Electrical Engineering at the California Institute of Technology. He received a B.Eng. and a M.Sc. in Electronic Engineering from Tsinghua University, China, in 2013 and 2016 respectively. His research interests include information and coding theory, with its applications in data storage systems. He is a recipient of the 2019 IEEE Jack Keil Wolf ISIT Student Paper Award.

**Jehoshua Bruck** (S'86-M'89-SM'93-F'01) is the Gordon and Betty Moore Professor of computation and neural systems and electrical engineering at the California Institute of Technology (Caltech). His current research interests include information theory and systems and the theory of computation in nature.

Dr. Bruck received the B.Sc. and M.Sc. degrees in electrical engineering from the Technion-Israel Institute of Technology, in 1982 and 1985, respectively, and the Ph.D. degree in electrical engineering from Stanford University, in 1989. His industrial and entrepreneurial experiences include working with IBM Research where he participated in the design and implementation of the first IBM parallel computer; cofounding and serving as Chairman of Rainfinity (acquired in 2005 by EMC), a spin-off company from Caltech that created the first virtualization solution for Network Attached Storage; as well as cofounding and serving as Chairman of XtremIO (acquired in 2012 by EMC), a start-up company that created the first scalable all-flash enterprise storage system.